

# PATENT ABSTRACTS OF JAPAN

(11)Publication number : **2002-091477**

(43)Date of publication of application : **27.03.2002**

(51)Int.Cl. **G10L 15/06**  
**G10L 15/14**  
**G10L 15/18**  
**G10L 15/28**

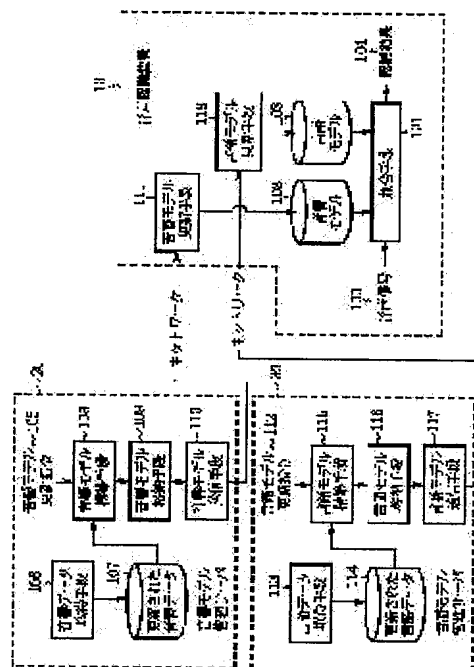
(21)Application number : **2000-280674** (71)Applicant : **MITSUBISHI ELECTRIC CORP**

(22)Date of filing : **14.09.2000** (72)Inventor : **OKATO YOHEI**  
**ISHII JUN**

**(54) VOICE RECOGNITION SYSTEM, VOICE RECOGNITION DEVICE, ACOUSTIC MODEL CONTROL SERVER, LANGUAGE MODEL CONTROL SERVER, VOICE RECOGNITION METHOD AND COMPUTER READABLE RECORDING MEDIUM WHICH RECORDS VOICE RECOGNITION PROGRAM**

(57)Abstract:

**PROBLEM TO BE SOLVED:** To update acoustic and language models and to improve the precision in recognition without exerting a large load onto a user.  
**SOLUTION:** An acoustic model control server 20 obtains updated acoustic data 107 and constructs an acoustic model. A language model control server 30 obtains updated language data 114 and constructs a language model. These models are respectively transmitted to a voice recognition device 10. In the device 10, an acoustic model updating means 111 updates an acoustic model 102 by the transmitted acoustic model. A language model updating means 118 updates a language model 103 using the transmitted language model.



## \* NOTICES \*

JPO and INPIT are not responsible for any damages caused by the use of this translation.

- 1.This document has been translated by computer. So the translation may not reflect the original precisely.
- 2.\*\*\* shows the word which can not be translated.
- 3.In the drawings, any words are not translated.

---

CLAIMS

---

## [Claim(s)]

[Claim 1]Voice recognition equipment which inputs an audio signal, performs speech recognition with reference to an acoustic model which asks for probability of an audio acoustical observed value series, and outputs a recognition result.

An acoustic model managing server which is connected with the above-mentioned voice recognition equipment via a network, acquires updated sound data, and builds the above-mentioned acoustic model.

It updates by an acoustic model to which it is the voice recognition system provided with the above, the above-mentioned acoustic model which the above-mentioned acoustic model managing server built was transmitted to the above-mentioned voice recognition equipment, and the above-mentioned acoustic model managing server transmitted an acoustic model which the above-mentioned voice recognition equipment refers to in the case of speech recognition.

[Claim 2]It corresponds to a specific condition directed by ID from which an acoustic model managing server acquired and acquired ID which specifies an acoustic model which voice recognition equipment refers to in the case of speech recognition, The voice recognition system according to claim 1 reading updated sound data, building an acoustic model depending on the above-mentioned specific condition, and transmitting to the above-mentioned voice recognition equipment.

[Claim 3]Voice recognition equipment which inputs an audio signal, performs speech recognition with reference to a language model which searches for appearing probability of a word sequence, and outputs a recognition result.

A language model managing server which is connected with the above-mentioned voice recognition equipment via a network, acquires updated language data, and builds the above-mentioned language model.

It updates by a language model to which it is the voice recognition system provided with the above, the above-mentioned language model which the above-mentioned language model managing server built was transmitted to the above-mentioned voice recognition equipment, and the above-mentioned language model managing server transmitted a language model which the above-mentioned voice recognition equipment refers to in the case of speech recognition.

[Claim 4]It corresponds to a specific condition directed by ID from which a language model managing server acquired and acquired ID which specifies a language model which voice recognition equipment refers to in the case of speech recognition, The voice recognition system according to claim 3 reading updated language data, building a language model depending on the above-mentioned specific condition, and transmitting to the above-mentioned voice recognition equipment.

[Claim 5]Language data in which a language model managing server read the above-mentioned user dictionary, and was updated via a network with reference to a user dictionary with which voice recognition equipment registered a word on the occasion of speech recognition, The voice recognition system according to claim 3 building a language model depending on the above-mentioned user dictionary with reference to the read above-mentioned user dictionary, and transmitting to the above-mentioned voice recognition equipment.

[Claim 6]The voice recognition system according to claim 3, wherein a language model managing server acquires a text which a user of voice recognition equipment used, builds a language model depending on the above-mentioned text with reference to updated language data and the acquired above-mentioned text and transmits to the above-mentioned voice recognition equipment.

[Claim 7]Voice recognition equipment which inputs an audio signal, performs speech recognition with reference to an acoustic model which asks for probability of an audio acoustical observed value series, and outputs a recognition result.

An acoustic model managing server which is connected with the above-mentioned voice recognition equipment via a network, and has an initial acoustic model before adaptation.

ID as which it is the voice recognition system provided with the above, and the above-mentioned voice recognition equipment specifies the above-mentioned acoustic model, Acquire voice data for adaptation from an inputted audio signal, and voice data acquired ID and for adaptation, Transmit to the above-mentioned acoustic model managing server via a network, and the above-mentioned acoustic model managing server, Using transmitted voice data for adaptation, carry out adaptation of the above-mentioned initial acoustic model, and match and store an adaptation finishing acoustic model in above-mentioned ID to which it was transmitted, and. In response to renewal instructions of an acoustic model from the outside, ID which specifies the above-mentioned acoustic model from the above-mentioned voice recognition equipment via a network is received, An adaptation finishing acoustic model corresponding to ID which received is chosen and read out of a stored adaptation finishing acoustic model, It updates by an adaptation finishing acoustic model to which it transmitted to the above-mentioned voice recognition equipment via a network, and the above-mentioned acoustic model managing server transmitted an acoustic model which the above-mentioned voice recognition equipment refers to in the case of speech recognition.

[Claim 8]An acoustic model which asks for probability of an audio acoustical observed value series.

A collation means which inputs an audio signal, performs speech recognition with reference to the above-mentioned acoustic model, and outputs a recognition result.

From an acoustic model managing server which is voice recognition equipment provided with the above, and was connected via a network. It had an acoustic model update means updated by an acoustic model which received an acoustic model built with updated sound data, and received an acoustic model which the above-mentioned collation means refers to in the case of speech recognition.

[Claim 9]An acoustic model update means from an acoustic model managing server connected via a network. The voice recognition equipment according to claim 8 updating by an acoustic model which received an acoustic model depending on a specific condition of an acoustic model which a collation means refers to in the case of speech recognition built with updated sound data, and received an acoustic model which the above-mentioned collation means refers to in the case of speech recognition.

[Claim 10]A language model which searches for appearing probability of a word sequence.

A collation means which inputs an audio signal, performs speech recognition with reference to the above-mentioned language model, and outputs a recognition result.

From a language model managing server which is voice recognition equipment provided with the above, and was connected via a network. It had a

language model update means updated by a language model which received a language model built with updated language data, and received a language model which the above-mentioned collation means refers to in the case of speech recognition.

[Claim 11] A language model update means from a language model managing server connected via a network. The voice recognition equipment according to claim 10 updating by a language model which received a language model depending on a specific condition of a language model which a collation means refers to in the case of speech recognition built with updated language data, and received a language model which the above-mentioned collation means refers to in the case of speech recognition.

[Claim 12] A collation means is provided with a user dictionary which registered a word referred to in the case of speech recognition, A language model update means from a language model managing server connected via a network. The voice recognition equipment according to claim 10 updating by a language model which received a language model depending on a user dictionary which was built with updated language data, and which the above-mentioned collation means refers to in the case of speech recognition, and received a language model which the above-mentioned collation means refers to in the case of speech recognition.

[Claim 13] A language model update means from a language model managing server connected via a network. The voice recognition equipment according to claim 10 updating by a language model which received a language model depending on a text which was built with updated language data, and which a user who performs speech recognition used, and received a language model which the above-mentioned collation means refers to in the case of speech recognition.

[Claim 14] Voice recognition equipment comprising:

An acoustic model which asks for probability of an audio acoustical observed value series.

A collation means which inputs an audio signal, performs speech recognition with reference to the above-mentioned acoustic model, and outputs a recognition result.

An acoustic model ID acquiring means which acquires ID which specifies the above-mentioned acoustic model.

Read ID which the above-mentioned acoustic model ID acquiring means acquired, and voice data for adaptation is acquired from an inputted audio signal, A voice acquisition means for adaptation which transmits voice data read ID and for acquired adaptation to an acoustic model managing server to which it was connected via a network, An acoustic model update means updated by an adaptation finishing acoustic model which received an adaptation finishing acoustic model by which adaptation was carried out with voice data for the above-mentioned adaptation corresponding to above-mentioned ID from the above-mentioned acoustic model managing server, and received an acoustic model which the above-mentioned collation means refers to in the case of speech recognition.

[Claim 15] An acoustic model managing server comprising:

A sound data acquisition means which acquires updated sound data.

An acoustic model construction means which builds an acoustic model which reads updated sound data which the above-mentioned sound data acquisition means acquired in response to renewal instructions of an acoustic model from the outside, and asks for probability of an audio acoustical observed value series.

An acoustic model transmitting means which transmits an acoustic model built by the above-mentioned acoustic model construction means to voice recognition equipment which performs speech recognition via a network.

[Claim 16] An acoustic model managing server comprising:

A sound data acquisition means which acquires updated sound data.

An updating acoustic model ID acquiring means which acquires ID which specifies an acoustic model which voice recognition equipment connected via a network refers to in response to renewal instructions of an acoustic model from the outside in the case of speech recognition.

A sound data reading means for the specification which reads updated sound data which the above-mentioned sound data acquisition means acquired corresponding to a specific condition directed by ID which the above-mentioned updating acoustic model ID acquiring means acquired.

An acoustic model construction means for the specification which builds an acoustic model depending on the above-mentioned specific condition with reference to updated sound data which the above-mentioned sound data reading means for specification read, An acoustic model transmitting means which transmits an acoustic model which the above-mentioned acoustic model construction means for specification built to the above-mentioned voice recognition equipment via a network.

[Claim 17] An acoustic model managing server comprising:

An initial acoustic model before adaptation which asks for probability of an audio acoustical observed value series.

Voice data for adaptation transmitted from voice recognition equipment connected via a network.

An acoustic model adaptation means for the above-mentioned voice recognition equipment to receive ID which specifies an acoustic model referred to in the case of speech recognition, to carry out adaptation of the above-mentioned initial acoustic model using voice data for the above-mentioned adaptation, to match an adaptation finishing acoustic model with above-mentioned ID which received, and to store in an adaptation finishing acoustic model storing means.

In response to renewal instructions of an acoustic model from the outside, above-mentioned ID is received from the above-mentioned voice recognition equipment via a network, An adaptation finishing acoustic model selecting means which chooses and reads an adaptation finishing acoustic model corresponding to ID which received from the above-mentioned adaptation finishing acoustic model storing means, An acoustic model transmitting means which transmits an adaptation finishing acoustic model which the above-mentioned adaptation finishing acoustic model selecting means read to the above-mentioned voice recognition equipment via a network.

[Claim 18] A language model managing server comprising:

A language data acquisition means which acquires updated language data.

A language model construction means which builds a language model which reads updated language data which the above-mentioned language data acquisition means acquired in response to renewal instructions of a language model from the outside, and searches for appearing probability of a word sequence.

A language model transmitting means which transmits a language model which the above-mentioned language model construction means built to voice recognition equipment which performs speech recognition via a network.

[Claim 19] A language model managing server comprising:

A language data acquisition means which acquires updated language data.

An updating language model ID acquiring means which acquires ID which specifies a language model which voice recognition equipment connected via a network refers to in response to renewal instructions of a language model from the outside in the case of speech recognition.

A language data reading means for the specification which reads updated language data which the above-mentioned language data acquisition means acquired corresponding to a specific condition directed by ID which the above-mentioned updating language model ID acquiring means acquired.

A language model construction means for the specification which builds a language model depending on the above-mentioned specific condition with reference to updated language data which the above-mentioned language data reading means for specification read, A language model transmitting

means which transmits a language model which the above-mentioned language model construction means for specification built to the above-mentioned voice recognition equipment via a network.

[Claim 20]A language model managing server comprising:

A language data acquisition means which acquires updated language data.

A user dictionary reading means which reads a user dictionary which voice recognition equipment connected via a network refers to in response to renewal instructions of a language model from the outside in the case of speech recognition.

A user dictionary dependence language model construction means which builds a language model depending on a user dictionary which read updated language data which the above-mentioned language data acquisition means acquired, and the above-mentioned user dictionary reading means read.

A language model transmitting means which transmits a language model which the above-mentioned user dictionary dependence language model construction means built to the above-mentioned voice recognition equipment via a network.

[Claim 21]A language model managing server comprising:

A language data acquisition means which acquires updated language data.

A user use text acquisition means which acquires a text which a user of voice recognition equipment connected via a network used in response to renewal instructions of a language model from the outside.

A user use text dependence language model construction means which builds a language model depending on a text which read updated language data which the above-mentioned language data acquisition means acquired, and the above-mentioned user use text acquisition means acquired.

A language model transmitting means which transmits a language model which the above-mentioned user use text dependence language model construction means built to the above-mentioned voice recognition equipment via a network.

[Claim 22]A voice recognition method which inputs an audio signal, performs speech recognition with reference to an acoustic model which asks for probability of an audio acoustical observed value series, and outputs a recognition result, comprising:

The 1st step that acquires updated sound data.

The 2nd step that reads updated sound data which was acquired at the 1st step of the above in response to renewal instructions of an acoustic model, and builds an acoustic model.

The 3rd step that transmits an acoustic model built at the 2nd step of the above via a network.

The 4th step updated by an acoustic model which received an acoustic model which transmitted at the 3rd step of the above, and received an acoustic model referred to in the case of the above-mentioned speech recognition.

[Claim 23]A voice recognition method which inputs an audio signal, performs speech recognition with reference to a language model which searches for appearing probability of a word sequence, and outputs a recognition result, comprising:

The 1st step that acquires updated language data.

The 2nd step that reads updated language data which was acquired at the 1st step of the above in response to renewal instructions of a language model, and builds a language model.

The 3rd step that transmits a language model built at the 2nd step of the above via a network.

The 4th step updated by a language model which received a language model which transmitted at the 3rd step of the above, and received a language model referred to in the case of the above-mentioned speech recognition.

[Claim 24]A voice recognition method which inputs an audio signal, performs speech recognition with reference to an acoustic model which asks for probability of an audio acoustical observed value series, and outputs a recognition result, comprising:

The 1st step that acquires updated sound data.

The 2nd step that acquires ID which specifies an acoustic model referred to in the case of speech recognition in response to renewal instructions of an acoustic model.

The 3rd step that reads updated sound data which was acquired at the 1st step of the above corresponding to a specific condition directed by ID acquired at the 2nd step of the above.

The 4th step that builds an acoustic model depending on the above-mentioned specific condition with reference to updated sound data which was read at the 3rd step of the above.

The 5th step that transmits an acoustic model built at the 4th step of the above via a network.

The 6th step updated by an acoustic model which received an acoustic model which transmitted at the 5th step of the above, and received an acoustic model referred to in the case of speech recognition.

[Claim 25]A voice recognition method which inputs an audio signal, performs speech recognition with reference to a language model which searches for appearing probability of a word sequence, and outputs a recognition result, comprising:

The 1st step that acquires updated language data.

The 2nd step that acquires ID which specifies a language model referred to in the case of speech recognition in response to renewal instructions of a language model.

The 3rd step that reads updated language data which was acquired at the 1st step of the above corresponding to a specific condition directed by ID acquired at the 2nd step of the above.

The 4th step that builds a language model depending on the above-mentioned specific condition with reference to updated language data which was read at the 3rd step of the above.

The 5th step that transmits a language model built at the 4th step of the above via a network.

The 6th step updated by a language model which received a language model which transmitted at the 5th step of the above, and received a language model referred to in the case of speech recognition.

[Claim 26]A voice recognition method which inputs an audio signal characterized by comprising the following, performs speech recognition with reference to a language model which searches for appearing probability of a word sequence, and a user dictionary which registered a word, and outputs a recognition result.

The 1st step that acquires updated language data.

The 2nd step that reads a user dictionary referred to in the case of speech recognition in response to renewal instructions of a language model.

The 3rd step that builds a language model depending on a user dictionary which read updated language data which was acquired at the 1st step of the above, and was read at the 2nd step of the above.

The 4th step that transmits a language model built at the 3rd step of the above via a network.

The 5th step updated by a language model which received a language model which transmitted at the 4th step of the above, and received a language model referred to in the case of speech recognition.

[Claim 27]A voice recognition method which inputs an audio signal, performs speech recognition with reference to a language model which searches for appearing probability of a word sequence, and outputs a recognition result, comprising:

The 1st step that acquires updated language data.

The 2nd step that acquires a text which a user who performs speech recognition used in response to renewal instructions of a language model.

The 3rd step that builds a language model depending on a text which read updated language data which was acquired at the 1st step of the above, and was acquired at the 2nd step of the above.

The 4th step that transmits a language model built at the 3rd step of the above via a network.

The 5th step updated by a language model which received a language model which transmitted at the 4th step of the above, and received a language model referred to in the case of speech recognition.

[Claim 28]A voice recognition method which inputs an audio signal, performs speech recognition with reference to an acoustic model which asks for probability of an audio acoustical observed value series, and outputs a recognition result, comprising:

The 1st step that acquires ID which specifies the above-mentioned acoustic model.

The 2nd step that transmits voice data ID which read ID acquired at the 1st step of the above, acquired voice data for adaptation from an inputted audio signal, and was read via a network, and for acquired adaptation.

The 3rd step that carries out adaptation of the initial acoustic model before adaptation, and matches and stores an adaptation finishing acoustic model in ID which transmitted at the 2nd step of the above using voice data for adaptation transmitted at the 2nd step of the above.

The 4th step that receives ID acquired at the 1st step of the above via a network, and chooses and reads an adaptation finishing acoustic model corresponding to ID which received in response to renewal instructions of an acoustic model out of an adaptation finishing acoustic model stored at the 3rd step of the above.

The 5th step that transmits an adaptation finishing acoustic model read at the 4th step of the above via a network.

The 6th step updated by an adaptation finishing acoustic model which received an adaptation finishing acoustic model which transmitted at the 5th step of the above, and received an acoustic model referred to in the case of speech recognition.

[Claim 29]Input an audio signal and speech recognition is performed with reference to an acoustic model which asks for probability of an audio acoustical observed value series, A sound data acquisition function which is the recording medium which recorded a voice recognition program which realizes a verification function which outputs a recognition result, and acquires updated sound data, An acoustic model construction function to read updated sound data which the above-mentioned sound data acquisition function acquired in response to renewal instructions of an acoustic model, and to build an acoustic model, An acoustic model transmitting function which transmits an acoustic model which the above-mentioned acoustic model construction function built via a network, A recording medium which recorded a voice recognition program which realizes an acoustic model update function updated by an acoustic model which received an acoustic model which the above-mentioned acoustic model transmitting function transmitted, and received an acoustic model which the above-mentioned verification function refers to in the case of speech recognition and in which computer reading is possible.

[Claim 30]Input an audio signal and speech recognition is performed with reference to a language model which searches for appearing probability of a word sequence, A language data acquisition function which is the recording medium which recorded a voice recognition program which realizes a verification function which outputs a recognition result, and acquires updated language data, A language model construction function to read updated language data which the above-mentioned language data acquisition function acquired in response to renewal instructions of a language model, and to build a language model, A language model transmitting function which transmits a language model which the above-mentioned language model construction function built via a network, A recording medium which recorded a voice recognition program which realizes a language model update function updated by a language model which received a language model which the above-mentioned language model transmitting function transmitted, and received a language model which the above-mentioned verification function refers to in the case of speech recognition and in which computer reading is possible.

[Claim 31]Input an audio signal and speech recognition is performed with reference to an acoustic model which asks for probability of an audio acoustical observed value series, A sound data acquisition function which is the recording medium which recorded a voice recognition program which realizes a verification function which outputs a recognition result, and acquires updated sound data, An updating acoustic model ID acquisition function which acquires ID which specifies the above-mentioned acoustic model in response to renewal instructions of an acoustic model, It corresponds to a specific condition directed by ID which the above-mentioned updating acoustic model ID acquisition function acquired, A sound data read-out function for the specification which reads updated sound data which the above-mentioned sound data acquisition function acquired, An acoustic model construction function for the specification which builds an acoustic model depending on the above-mentioned specific condition with reference to updated sound data which the above-mentioned sound data read-out function for specification read, An acoustic model transmitting function which transmits an acoustic model which the above-mentioned acoustic model construction function for specification built via a network, A recording medium which recorded a voice recognition program which realizes an acoustic model update function updated by an acoustic model which received an acoustic model which the above-mentioned acoustic model transmitting function transmitted, and received an acoustic model which the above-mentioned verification function refers to in the case of speech recognition and in which computer reading is possible.

[Claim 32]Input an audio signal and speech recognition is performed with reference to a language model which searches for appearing probability of a word sequence, A language data acquisition function which is the recording medium which recorded a voice recognition program which realizes a verification function which outputs a recognition result, and acquires updated language data, An updating language model ID acquisition function which acquires ID which specifies the above-mentioned language model in response to renewal instructions of a language model, It corresponds to a specific condition directed by ID which the above-mentioned updating language model ID acquisition function acquired, A language data read-out function for the specification which reads updated language data which the above-mentioned language data acquisition function acquired, A language model construction function for the specification which builds a language model depending on the above-mentioned specific condition with reference to updated language data which the above-mentioned language data read-out function for specification read, A language model transmitting function which transmits a language model which the above-mentioned language model construction function for specification built via a network, A recording medium which recorded a voice recognition program which realizes a language model update function updated by a language model which received a language model which the above-mentioned language model transmitting function transmitted, and received a language model which the above-mentioned verification function refers to in the case of speech recognition and in which computer reading is possible.

[Claim 33]Input an audio signal and speech recognition is performed with reference to a language model which searches for appearing probability of a word sequence, and a user dictionary which registered a word, A language data acquisition function which is the recording medium which recorded a voice recognition program which realizes a verification function which outputs a recognition result, and acquires updated language data, A user dictionary read-out function which reads the above-mentioned user dictionary in response to renewal instructions of a language model, A user dictionary dependence language model construction function to build a language model depending on a user dictionary which read updated language data which the above-mentioned language data acquisition function acquired, and the above-mentioned user dictionary read-out function read, A language model transmitting function which transmits a language model which the above-mentioned user dictionary dependence language model construction function built via a network, A recording medium which recorded a voice recognition program which realizes a language model update function updated by a language model which received a language model which the above-mentioned language model transmitting function transmitted, and received a language model which the above-mentioned verification function refers to in the case of speech recognition and in which computer

reading is possible.

[Claim 34] Input an audio signal and speech recognition is performed with reference to a language model which searches for appearing probability of a word sequence, A language data acquisition function which is the recording medium which recorded a voice recognition program which realizes a verification function which outputs a recognition result, and acquires updated language data, A user use text acquisition function which acquires a text which a user who performs speech recognition used in response to renewal instructions of a language model, A user use text dependence language model construction function to build a language model depending on a text which read updated language data which the above-mentioned language data acquisition function acquired, and the above-mentioned user use text acquisition function acquired, A language model transmitting function which transmits a language model which the above-mentioned user use text dependence language model construction function built via a network, A recording medium which recorded a voice recognition program which realizes a language model update function updated by a language model which received a language model which the above-mentioned language model transmitting function transmitted, and received a language model which the above-mentioned verification function refers to in the case of speech recognition and in which computer reading is possible.

[Claim 35] Input an audio signal and speech recognition is performed with reference to an acoustic model which asks for probability of an audio acoustical observed value series, An acoustic model ID acquisition function which is the recording medium which recorded a voice recognition program which realizes a verification function which outputs a recognition result, and acquires ID which specifies the above-mentioned acoustic model, Read ID which the above-mentioned acoustic model ID acquisition function acquired, acquire voice data for adaptation from an inputted audio signal, and via a network, A voice acquisition function for adaptation which transmits voice data read ID and for acquired adaptation, Voice data for adaptation which the above-mentioned voice acquisition function for adaptation transmitted is used, An acoustic model adaptation function to carry out adaptation of the initial acoustic model before adaptation, and to match and store an adaptation finishing acoustic model in ID which the above-mentioned voice acquisition function for adaptation transmitted, In response to renewal instructions of an acoustic model, ID which the above-mentioned acoustic model ID acquisition function acquired via a network is received, An adaptation finishing acoustic model which an adaptation finishing acoustic model function preselection capability which chooses and reads an adaptation finishing acoustic model corresponding to ID which received out of an adaptation finishing acoustic model which the above-mentioned acoustic model adaptation function stored, and the above-mentioned adaptation finishing acoustic model function preselection capability read, An adaptation finishing acoustic model which an acoustic model transmitting function which transmits via a network, and the above-mentioned acoustic model transmitting function transmitted is received, A recording medium which recorded a voice recognition program which realizes an acoustic model update function in which the above-mentioned verification function updates an acoustic model referred to in the case of speech recognition by an adaptation finishing acoustic model which received and in which computer reading is possible.

---

[Translation done.]

## \* NOTICES \*

JPO and INPIT are not responsible for any damages caused by the use of this translation.

- 1.This document has been translated by computer. So the translation may not reflect the original precisely.
- 2.\*\*\* shows the word which can not be translated.
- 3.In the drawings, any words are not translated.

## DETAILED DESCRIPTION

[Detailed Description of the Invention]

[0001]

[Field of the Invention]This invention the acoustic model and language model which are referred to in the case of speech recognition so that a high recognition rate may be acquired, It is related with the recording medium which recorded the voice recognition system, the voice recognition equipment, the acoustic model managing server, the language model managing server, voice recognition method, and voice recognition program which are updated in the newest state via a network and in which computer reading is possible.

[0002]

[Description of the Prior Art]In speech recognition, after usually changing the sound digitized and inputted into the time series of the vector which expresses the audio acoustical feature well using the signal-processing technique, collation processing with an audio model (an acoustic model, a language model) is performed.

[0003]Collation processing is a problem which asks for word sequence  $W = [w_1 w_2 \dots w_k]$  ( $k$  is the number of words) uttered from sound feature vector time series  $A = [a_1 a_2 \dots a_n]$  which consists of  $n$  time frames. In order to presume a word sequence as for which recognition precision becomes the highest, the appearing probability  $P(W|A)$  should just ask for recognized word sequence  $W^*$  used as the maximum. Namely,  $W^* = \arg\max W P(W|A)$  (1)

However, it is usually difficult to ask for  $P(W|A)$  directly. Then,  $P(W|A)$  is rewritten like (2) types using the Bayes theorem.

$$P(W|A) = P(W)P(A|W)/P(A) \quad (2)$$

[0004]Here, when asking for  $W$  which maximizes the left side, since  $P(A)$  which is a right-hand side denominator does not affect  $W$  which becomes a recognition candidate, it should just ask for  $W$  which maximizes a right-hand side numerator. That is, it becomes like (3) types.

$$W^* = \arg\max W P(W)P(A|W) \quad (3)$$

Here, the stochastic model which gives  $P(W)$  is called a language model, and the stochastic model which gives  $P(A|W)$  is called an acoustic model.

[0005]These typical modeling methods in speech recognition are the methods of expressing an acoustic model by Hidden Markov Model, and expressing a language model by  $n-1$ -fold Markov process of the word called  $n$  gram.

[0006]The details of these methods, for example "Basic [ of speech recognition ] (top, under)" L.RABINER, B.H.JUANG, the Furui supervision of translation, November, 1995, NTT Advanced Technology (it is hereafter considered as the literature 1), It is describing in Shokodo (it is hereafter considered as the literature 3) in the "stochastic language model" Kitasato Institute 2, University of Tokyo Press (it is hereafter considered as the literature 2), "digital signal processing of sound and sound information" Kiyohiro Kano, Tetsu Nakamura, Shiro Ise collaboration, and November, 1997.

[0007]In these methods, the parameter which constitutes a model is statistically presumed from a lot of data. In construction of an acoustic model, voice data, such as a word from many speakers and a sentence, is collected beforehand, and it is presumed that the index well connected with recognition precision or recognition precision using the statistical method improves. For example, the parameter of the Hidden Markov Model which constitutes an acoustic model is presumed using a BAUMU Welch algorithm so that the likelihood to which an acoustic model outputs learned data may serve as the maximum. The estimation method of the acoustic model is described in detail in literature 1 beaming.

[0008]Similarly, in construction of a language model, the probability that the word which constitutes each utterance and utterance will appear from texts, such as a newspaper and the opening sentence of conversation, according to the structure of a language model is calculated. For example, in  $n$  gram language model, when it sets with  $n = 2$  (called a bigram language model),  $P(W)$  is approximated like (4) types. The parameter of  $n$  gram language model is presumed from the frequency of adjoining  $n$  word in the text data for study. The estimation method of the language model is described in detail in the literature 2.

$$P(w_1 \dots w_k)$$

$$= P(w_1)P(w_2|w_1) \dots P(w_k|w_1 \dots w_{k-1})$$

$$= P(w_1)P(w_2|w_1) \dots P(w_k|w_{k-1}) \quad (4)$$

[0009]Thus, the language model which acquires high recognition precision compared with the method of not using a statistics value can be built by presuming the appearing probability of each word statistically using an extensive text. Since a text is not written with a space between words in Japanese, the definition of a word is ambiguous, but in the text, each unit which divided the text by a certain compatible means is defined as a word. This word is the division of a text based on linguistic units, such as a character, a morpheme, a clause, or an entropy standard, such combination, etc., for example.

[0010]Drawing 12 is a block diagram showing the composition of the conventional voice recognition equipment indicated by the above-mentioned literature 1. In drawing 12, the acoustic model referred to when the collation means 101 carries out speech recognition of the collation means which 101 inputs and carries out speech recognition of the audio signal 100, and outputs the recognition result 104, and 102, and 103 are language models referred to when the collation means 101 carries out speech recognition.

[0011]Next, operation is explained. The collation means 101 inputs a user's audio signal 100, performs speech recognition with reference to the acoustic model 102 and the language model 103, and outputs the recognition result 104. The acoustic model 102 expresses the mapping relation of the time series of the sound feature vector produced by carrying out signal processing of the voice waveform of a user's inputted audio signal 100, and the minimum symbol information that the voice recognition equipment expressed with a phoneme etc. treats, for example. The language model 103 describes the appearance information of a word to be a correspondence relation with recognition units longer than the word etc. which are expressed with the combination of the symbol mapped by the acoustic model 102. The acoustic model 102 asks for the probability that the model of a certain symbol will output a vector time series.

That is, it asks for the probability of an audio acoustical observed value series, and the language model 103 searches for the appearing probability of a certain word sequence.

[0012]Drawing 13 is a flow chart which shows the procedure of the conventional speech recognition processing in which input the audio signal 100 and the recognition result 104 is obtained. In step ST1301, the A/D conversion of the inputted audio signal 100 is carried out, and it turns into a digital signal. In step ST1302, the digitized audio signal sets a suitable interval, and signal processing is carried out, and it is changed into the time series of the sound feature vector which expresses audio character well.

[0013]In step ST1303, a sound feature vector is compared with the acoustic model 102 by sound collation processing, and the probability which outputs the time series of a sound feature vector about each recognition candidate is called for. In step ST1304, further, each recognition candidate is compared with the language model 103 by language collation processing, and the output probability of a word sequence can multiply by him by it. Finally,



in step ST1305, the most suitable candidate is chosen from each recognition candidate, and the recognition result 104 is obtained. Usually, the most suitable recognition result is the recognition candidate made the highest [ probability ] by the above-mentioned collation.

[0014]Only by the acoustic model and language model which were constituted by the described method. In the voice recognition equipment which a user can customize by the case where sufficient performance is not attained. In order to make the sound and word which are hard to be recognized recognize better, recognition precision may be able to be raised by making a user do adaptation of the acoustic model, or adding a recognition object word to a user dictionary.

[0015]First, the case where adaptation of the acoustic model is carried out is explained. Drawing 14 is a block diagram showing the composition of conventional voice recognition equipment provided with an acoustic model adaptation means indicated by the above-mentioned literature 1 to carry out adaptation of the acoustic model to the audio signal 100. In drawing 14, a different portion from drawing 12 in order to adapt an acoustic model to the audio signal 100 to input, It is having the initial acoustic model 1003, the acoustic model adaptation means 1004, and the adaptation finishing acoustic model 1401 and that the audio signal 100 is chosen by the acoustic model adaptation means 1004 and the collation means 101.

[0016]Next, operation of the voice recognition equipment shown in drawing 14 is explained. For example using a maximum-a-posteriori-probability-estimation method, the acoustic model adaptation means 1004 carries out [ sound / for adaptation ] adaptation of the initial acoustic model 1003, and obtains the adaptation finishing acoustic model 1401 from the sounds for adaptation (audio signal 100) and the initial acoustic models 1003 collected before actual recognition. The adaptation method of the acoustic model is shown in Chapter 7 of the literature 3. The collation means 101 inputs the audio signal 100, performs speech recognition with reference to the adaptation finishing acoustic model 1401 and the language model 103, and outputs the recognition result 104.

[0017]Next, the case where a user registers a recognition object word into a dictionary is explained. Drawing 15 is a block diagram showing the composition of the conventional voice recognition equipment to which the user dictionary was added indicated by the above-mentioned literature 1. In drawing 15, a different portion from drawing 12 is having the user dictionary 601.

[0018]Drawing 16 is a figure showing the example of composition of the user dictionary 601. The user dictionary 601 is a meeting of the word which the user registered in order to make the word which is not recognized and the word which is hard to be recognized recognize better, and they are a notation of a word, and the list of the words which consist of reading.

[0019]Next, operation of the voice recognition equipment shown in drawing 15 is explained. Drawing 17 is a flow chart which shows the procedure of the conventional speech recognition processing in which the user dictionary 601 was added. In drawing 17, in the language collation processing of step ST1704, a different point from drawing 13 is that the word registered into the user dictionary 601 is added to a recognition object word, in order to refer to the language model 103 and the user dictionary 601.

[0020]The word registered into the user dictionary 601 enables arbitrary word sequences and connection according to suitable conjunctive probability. For example, in the bigram language model determined only with one word which precedes the appearance conditions of a word, Constant value is given to the probability  $P(\text{wuser}|\text{wi})$  and  $P$  of the user dictionary registration word "wuser" inserted into the arbitrary words "wi and wj" ( $\text{wjlwuser}$ ), and reallocation of the probability value is carried out so that the probability of the whole language model may be set to 1. As a result, the user can obtain the recognition result in which the registered word was contained.

[0021]However, in the composition of the voice recognition equipment shown in drawing 14 and drawing 15. In order to raise recognition precision, by creating the audio signal 100 and the adaptation finishing acoustic model 1401, or registering a word into the user dictionary 601, a user needs to customize the acoustic model 102 and the language model 103 himself, and will force a big burden upon a user.

[0022]Since there are a word which had not appeared, and direction for use which was not taken into consideration at the language model construction time even if it is a case where a word is registered to the user dictionary 601, the appearing probability given to those words may be unsuitable. This may become a cause and recognition precision may fall.

[0023]The acoustic model 102 and the language model 103 which were customized as mentioned above are referred to only from the specific collation means 101. For this reason, the other collation means's 101 use of this acoustic model 102 and language model 103 that were customized will reduce recognition precision.

[0024]For this reason, in language processing systems, such as Chinese character kana conversion and a machine translation. The system which eased the user's burden is proposed, for example like JP,10-260960,A by having a function which updates a dictionary automatically via a network about renewal of the user dictionary 601 among the problems shown above. However, in the above-mentioned gazette, it cannot take into consideration treating the model for pattern recognition like a sound, and it cannot be applied to the model about pattern information like the acoustic model 102.

[0025]Also in renewal of the language model 103, if a user's number of registered words increases, since [ which gushes and carries out and becomes easy to produce an error ] the short word which is not suitable was inserted, the fall of recognition precision will occur easily. Since it does not exist when the word which carried out additional registration to the user dictionary 601 builds the language model 103, or the operating environment is changing, the fall of this recognition precision. For example, it is because the appearing probability which is not suitable for the appearing probability  $P(\text{wuser}|\text{wi})$  and  $P$  of a user registration word ( $\text{wjlwuser}$ ), etc. may be given. As a result, it becomes easy to obtain an unsuitable recognition result, and recognition precision may fall. In order to prevent this, it is necessary to set up the appearing probability of a word appropriately but, and generally it is difficult for the user itself to give an appropriate value.

[0026]

[Problem(s) to be Solved by the Invention]When customizing an acoustic model and a language model in order to raise recognition precision since conventional voice recognition equipment was constituted as mentioned above, the technical problem that a big burden was applied to a user occurred.

[0027]When using except the voice recognition equipment customized for every user, the technical problem that recognition precision fell occurred.

[0028]The technical problem that an acoustic model could not be updated automatically by customization through a network occurred.

[0029]When the registration to a user dictionary increased, the technical problem that recognition precision fell easily occurred.

[0030]By the server side which was made in order that this invention might solve the above technical problems, and was connected to the network. Acquire the newest sound data or language data, and build the acoustic model or language model in the newest state, or, By building the acoustic model or language model corresponding to a user, and updating users' acoustic model and language model via a network. It aims at obtaining the recording medium which recorded the voice recognition system, the voice recognition equipment, the acoustic model managing server, the language model managing server, voice recognition method, and voice recognition program which raise recognition precision without applying a big burden to a user and in which computer reading is possible.

[0031]In all voice recognition equipments by connecting via a network. It aims at obtaining the recording medium which recorded the voice recognition system, the voice recognition equipment, the acoustic model managing server, the language model managing server, voice recognition method, and voice recognition program which can use the customized acoustic model or language model and in which computer reading is possible.

[0032]By using the dictionary and text which are obtained from a user, and the text collected semi-automatically. It aims at obtaining the recording medium which recorded the voice recognition system, the voice recognition equipment, the acoustic model managing server, the language model managing server, voice recognition method, and voice recognition program to which recognition precision cannot fall easily and in which computer reading is possible, even when a user dictionary becomes large.

[0033]

[Means for Solving the Problem]Voice recognition equipment which a voice recognition system concerning this invention inputs an audio signal, performs speech recognition with reference to an acoustic model which asks for probability of an audio acoustical observed value series, and outputs a recognition result. In a thing provided with an acoustic model managing server which is connected with the above-mentioned voice recognition equipment via a network, acquires updated sound data, and builds the above-mentioned acoustic model. The above-mentioned acoustic model which



the above-mentioned acoustic model managing server built is transmitted to the above-mentioned voice recognition equipment, and an acoustic model which the above-mentioned voice recognition equipment refers to in the case of speech recognition is updated by an acoustic model which the above-mentioned acoustic model managing server transmitted.

[0034]A voice recognition system concerning this invention corresponds to a specific condition directed by ID from which an acoustic model managing server acquired and acquired ID which specifies an acoustic model which voice recognition equipment refers to in the case of speech recognition, Updated sound data is read, an acoustic model depending on the above-mentioned specific condition is built, and it transmits to the above-mentioned voice recognition equipment.

[0035]Voice recognition equipment which a voice recognition system concerning this invention inputs an audio signal, performs speech recognition with reference to a language model which searches for appearing probability of a word sequence, and outputs a recognition result, In a thing provided with a language model managing server which is connected with the above-mentioned voice recognition equipment via a network, acquires updated language data, and builds the above-mentioned language model, The above-mentioned language model which the above-mentioned language model managing server built is transmitted to the above-mentioned voice recognition equipment, and a language model which the above-mentioned voice recognition equipment refers to in the case of speech recognition is updated by a language model which the above-mentioned language model managing server transmitted.

[0036]A voice recognition system concerning this invention corresponds to a specific condition directed by ID from which a language model managing server acquired and acquired ID which specifies a language model which voice recognition equipment refers to in the case of speech recognition, Updated language data is read, a language model depending on the above-mentioned specific condition is built, and it transmits to the above-mentioned voice recognition equipment.

[0037]Language data in which a language model managing server read the above-mentioned user dictionary via a network, and a voice recognition system concerning this invention was updated with reference to a user dictionary with which voice recognition equipment registered a word on the occasion of speech recognition, With reference to the read above-mentioned user dictionary, a language model depending on the above-mentioned user dictionary is built, and it transmits to the above-mentioned voice recognition equipment.

[0038]A language model managing server acquires a text which a user of voice recognition equipment used, and a voice recognition system concerning this invention builds a language model depending on the above-mentioned text with reference to updated language data and the acquired above-mentioned text, and transmits to the above-mentioned voice recognition equipment.

[0039]Voice recognition equipment which a voice recognition system concerning this invention inputs an audio signal, performs speech recognition with reference to an acoustic model which asks for probability of an audio acoustical observed value series, and outputs a recognition result, In a thing provided with an acoustic model managing server which is connected with the above-mentioned voice recognition equipment via a network, and has an initial acoustic model before adaptation, The above-mentioned voice recognition equipment acquires voice data for adaptation from an audio signal inputted as ID which specifies the above-mentioned acoustic model, Voice data acquired ID and for adaptation is transmitted to the above-mentioned acoustic model managing server via a network, Using transmitted voice data for adaptation, the above-mentioned acoustic model managing server carries out adaptation of the above-mentioned initial acoustic model, and matches and stores an adaptation finishing acoustic model in above-mentioned ID to which it was transmitted, and. In response to renewal instructions of an acoustic model from the outside, ID which specifies the above-mentioned acoustic model from the above-mentioned voice recognition equipment via a network is received, Choose and read an adaptation finishing acoustic model corresponding to ID which received out of a stored adaptation finishing acoustic model, transmit to the above-mentioned voice recognition equipment via a network, and the above-mentioned voice recognition equipment an acoustic model referred to in the case of speech recognition, It updates by an adaptation finishing acoustic model which the above-mentioned acoustic model managing server transmitted.

[0040]In a thing provided with a collation means which voice recognition equipment concerning this invention inputs an audio signal as an acoustic model which asks for probability of an audio acoustical observed value series, performs speech recognition with reference to the above-mentioned acoustic model, and outputs a recognition result, An acoustic model built with updated sound data from an acoustic model managing server connected via a network is received, and it has an acoustic model update means which updates an acoustic model which the above-mentioned collation means refers to in the case of speech recognition by an acoustic model which received.

[0041]An acoustic model update means voice recognition equipment concerning this invention from an acoustic model managing server connected via a network. An acoustic model depending on a specific condition of an acoustic model which a collation means refers to in the case of speech recognition built with updated sound data is received, and the above-mentioned collation means updates an acoustic model referred to in the case of speech recognition by an acoustic model which received.

[0042]In a thing provided with a collation means which voice recognition equipment concerning this invention inputs an audio signal as a language model which searches for appearing probability of a word sequence, performs speech recognition with reference to the above-mentioned language model, and outputs a recognition result, A language model built with updated language data from a language model managing server connected via a network is received, and it has a language model update means which updates a language model which the above-mentioned collation means refers to in the case of speech recognition by a language model which received.

[0043]A language model update means voice recognition equipment concerning this invention from a language model managing server connected via a network. A language model depending on a specific condition of a language model which a collation means refers to in the case of speech recognition built with updated language data is received, and the above-mentioned collation means updates a language model referred to in the case of speech recognition by a language model which received.

[0044]Voice recognition equipment concerning this invention is provided with a user dictionary which registered a word which a collation means refers to in the case of speech recognition, A language model update means from a language model managing server connected via a network. A language model depending on a user dictionary which was built with updated language data and which the above-mentioned collation means refers to in the case of speech recognition is received, and the above-mentioned collation means updates a language model referred to in the case of speech recognition by a language model which received.

[0045]A language model update means voice recognition equipment concerning this invention from a language model managing server connected via a network. A language model depending on a text which was built with updated language data and which a user who performs speech recognition used is received, and the above-mentioned collation means updates a language model referred to in the case of speech recognition by a language model which received.

[0046]Voice recognition equipment of this invention is characterized by comprising:

An acoustic model which asks for probability of an audio acoustical observed value series.

A collation means which inputs an audio signal, performs speech recognition with reference to the above-mentioned acoustic model, and outputs a recognition result.

An acoustic model ID acquiring means which acquires ID which specifies the above-mentioned acoustic model.

Read ID which the above-mentioned acoustic model ID acquiring means acquired, and voice data for adaptation is acquired from an inputted audio signal, A voice acquisition means for adaptation which transmits voice data read ID and for acquired adaptation to an acoustic model managing server to which it was connected via a network, An acoustic model update means updated by an adaptation finishing acoustic model which received an adaptation finishing acoustic model by which adaptation was carried out with voice data for the above-mentioned adaptation corresponding to above-mentioned ID from the above-mentioned acoustic model managing server, and received an acoustic model which the above-mentioned collation means refers to in the case of speech recognition.

[0047]An acoustic model managing server of this invention is characterized by comprising:

A sound data acquisition means which acquires updated sound data.

An acoustic model construction means which builds an acoustic model which reads updated sound data which the above-mentioned sound data acquisition means acquired in response to renewal instructions of an acoustic model from the outside, and asks for probability of an audio acoustical observed value series.

An acoustic model transmitting means which transmits an acoustic model built by the above-mentioned acoustic model construction means to voice recognition equipment which performs speech recognition via a network.

[0048]An acoustic model managing server of this invention is characterized by comprising:

A sound data acquisition means which acquires updated sound data.

An updating acoustic model ID acquiring means which acquires ID which specifies an acoustic model which voice recognition equipment connected via a network refers to in response to renewal instructions of an acoustic model from the outside in the case of speech recognition.

A sound data reading means for the specification which reads updated sound data which the above-mentioned sound data acquisition means acquired corresponding to a specific condition directed by ID which the above-mentioned updating acoustic model ID acquiring means acquired.

An acoustic model construction means for the specification which builds an acoustic model depending on the above-mentioned specific condition with reference to updated sound data which the above-mentioned sound data reading means for specification read, An acoustic model transmitting means which transmits an acoustic model which the above-mentioned acoustic model construction means for specification built to the above-mentioned voice recognition equipment via a network.

[0049]An acoustic model managing server of this invention is characterized by comprising:

An initial acoustic model before adaptation which asks for probability of an audio acoustical observed value series.

Voice data for adaptation transmitted from voice recognition equipment connected via a network.

An acoustic model adaptation means for the above-mentioned voice recognition equipment to receive ID which specifies an acoustic model referred to in the case of speech recognition, to carry out adaptation of the above-mentioned initial acoustic model using voice data for the above-mentioned adaptation, to match an adaptation finishing acoustic model with above-mentioned ID which received, and to store in an adaptation finishing acoustic model storing means.

In response to renewal instructions of an acoustic model from the outside, above-mentioned ID is received from the above-mentioned voice recognition equipment via a network, An adaptation finishing acoustic model selecting means which chooses and reads an adaptation finishing acoustic model corresponding to ID which received from the above-mentioned adaptation finishing acoustic model storing means, An acoustic model transmitting means which transmits an adaptation finishing acoustic model which the above-mentioned adaptation finishing acoustic model selecting means read to the above-mentioned voice recognition equipment via a network.

[0050]A language model managing server of this invention is characterized by comprising:

A language data acquisition means which acquires updated language data.

A language model construction means which builds a language model which reads updated language data which the above-mentioned language data acquisition means acquired in response to renewal instructions of a language model from the outside, and searches for appearing probability of a word sequence.

A language model transmitting means which transmits a language model which the above-mentioned language model construction means built to voice recognition equipment which performs speech recognition via a network.

[0051]A language model managing server of this invention is characterized by comprising:

A language data acquisition means which acquires updated language data.

An updating language model ID acquiring means which acquires ID which specifies a language model which voice recognition equipment connected via a network refers to in response to renewal instructions of a language model from the outside in the case of speech recognition.

A language data reading means for the specification which reads updated language data which the above-mentioned language data acquisition means acquired corresponding to a specific condition directed by ID which the above-mentioned updating language model ID acquiring means acquired.

A language model construction means for the specification which builds a language model depending on the above-mentioned specific condition with reference to updated language data which the above-mentioned language data reading means for specification read, A language model transmitting means which transmits a language model which the above-mentioned language model construction means for specification built to the above-mentioned voice recognition equipment via a network.

[0052]A language model managing server of this invention is characterized by comprising:

A language data acquisition means which acquires updated language data.

A user dictionary reading means which reads a user dictionary which voice recognition equipment connected via a network refers to in response to renewal instructions of a language model from the outside in the case of speech recognition.

A user dictionary dependence language model construction means which builds a language model depending on a user dictionary which read updated language data which the above-mentioned language data acquisition means acquired, and the above-mentioned user dictionary reading means read.

A language model transmitting means which transmits a language model which the above-mentioned user dictionary dependence language model construction means built to the above-mentioned voice recognition equipment via a network.

[0053]A language model managing server of this invention is characterized by comprising:

A language data acquisition means which acquires updated language data.

A user use text acquisition means which acquires a text which a user of voice recognition equipment connected via a network used in response to renewal instructions of a language model from the outside.

A user use text dependence language model construction means which builds a language model depending on a text which read updated language data which the above-mentioned language data acquisition means acquired, and the above-mentioned user use text acquisition means acquired.

A language model transmitting means which transmits a language model which the above-mentioned user use text dependence language model construction means built to the above-mentioned voice recognition equipment via a network.

[0054]A voice recognition method of this invention inputs an audio signal, performs speech recognition with reference to an acoustic model which asks for probability of an audio acoustical observed value series, and is characterized by that what outputs a recognition result comprises the following.

The 1st step that acquires updated sound data.

The 2nd step that reads updated sound data which was acquired at the 1st step of the above in response to renewal instructions of an acoustic model, and builds an acoustic model.

The 3rd step that transmits an acoustic model built at the 2nd step of the above via a network.

The 4th step updated by an acoustic model which received an acoustic model which transmitted at the 3rd step of the above, and received an acoustic

model referred to in the case of the above-mentioned speech recognition.

[0055]A voice recognition method of this invention inputs an audio signal, performs speech recognition with reference to a language model which searches for appearing probability of a word sequence, and is characterized by that what outputs a recognition result comprises the following.

The 1st step that acquires updated language data.

The 2nd step that reads updated language data which was acquired at the 1st step of the above in response to renewal instructions of a language model, and builds a language model.

The 3rd step that transmits a language model built at the 2nd step of the above via a network.

The 4th step updated by a language model which received a language model which transmitted at the 3rd step of the above, and received a language model referred to in the case of the above-mentioned speech recognition.

[0056]A voice recognition method of this invention inputs an audio signal, performs speech recognition with reference to an acoustic model which asks for probability of an audio acoustical observed value series, and is characterized by that what outputs a recognition result comprises the following.

The 1st step that acquires updated sound data.

The 2nd step that acquires ID which specifies an acoustic model referred to in the case of speech recognition in response to renewal instructions of an acoustic model.

The 3rd step that reads updated sound data which was acquired at the 1st step of the above corresponding to a specific condition directed by ID acquired at the 2nd step of the above.

The 4th step that builds an acoustic model depending on the above-mentioned specific condition with reference to updated sound data which was read at the 3rd step of the above, The 6th step updated by an acoustic model which received an acoustic model which transmitted an acoustic model built at the 4th step of the above at the 5th step that transmits via a network, and the 5th step of the above, and received an acoustic model referred to in the case of speech recognition.

[0057]A voice recognition method of this invention inputs an audio signal, performs speech recognition with reference to a language model which searches for appearing probability of a word sequence, and is characterized by that what outputs a recognition result comprises the following.

The 1st step that acquires updated language data.

The 2nd step that acquires ID which specifies a language model referred to in the case of speech recognition in response to renewal instructions of a language model.

The 3rd step that reads updated language data which was acquired at the 1st step of the above corresponding to a specific condition directed by ID acquired at the 2nd step of the above.

The 4th step that builds a language model depending on the above-mentioned specific condition with reference to updated language data which was read at the 3rd step of the above, The 6th step updated by a language model which received a language model which transmitted a language model built at the 4th step of the above at the 5th step that transmits via a network, and the 5th step of the above, and received a language model referred to in the case of speech recognition.

[0058]A voice recognition method of this invention inputs an audio signal, performs speech recognition with reference to a language model which searches for appearing probability of a word sequence, and a user dictionary which registered a word, and is characterized by that what outputs a recognition result comprises the following.

The 1st step that acquires updated language data.

The 2nd step that reads a user dictionary referred to in the case of speech recognition in response to renewal instructions of a language model.

The 3rd step that builds a language model depending on a user dictionary which read updated language data which was acquired at the 1st step of the above, and was read at the 2nd step of the above.

The 5th step updated by a language model which received a language model which transmitted a language model built at the 3rd step of the above at the 4th step that transmits via a network, and the 4th step of the above, and received a language model referred to in the case of speech recognition.

[0059]A voice recognition method of this invention inputs an audio signal, performs speech recognition with reference to a language model which searches for appearing probability of a word sequence, and is characterized by that what outputs a recognition result comprises the following.

The 1st step that acquires updated language data.

The 2nd step that acquires a text which a user who performs speech recognition used in response to renewal instructions of a language model.

The 3rd step that builds a language model depending on a text which read updated language data which was acquired at the 1st step of the above, and was acquired at the 2nd step of the above.

The 5th step updated by a language model which received a language model which transmitted a language model built at the 3rd step of the above at the 4th step that transmits via a network, and the 4th step of the above, and received a language model referred to in the case of speech recognition.

[0060]A voice recognition method of this invention inputs an audio signal, performs speech recognition with reference to an acoustic model which asks for probability of an audio acoustical observed value series, and is characterized by that what outputs a recognition result comprises the following.

The 1st step that acquires ID which specifies the above-mentioned acoustic model.

The 2nd step that transmits voice data ID which read ID acquired at the 1st step of the above, acquired voice data for adaptation from an inputted audio signal, and was read via a network, and for acquired adaptation.

The 3rd step that carries out adaptation of the initial acoustic model before adaptation, and matches and stores an adaptation finishing acoustic model in ID which transmitted at the 2nd step of the above using voice data for adaptation transmitted at the 2nd step of the above.

In response to renewal instructions of an acoustic model, ID acquired at the 1st step of the above via a network is received, The 4th step that chooses and reads an adaptation finishing acoustic model corresponding to ID which received out of an adaptation finishing acoustic model stored at the 3rd step of the above, The 5th step that transmits an adaptation finishing acoustic model read at the 4th step of the above via a network, The 6th step updated by an adaptation finishing acoustic model which received an adaptation finishing acoustic model which transmitted at the 5th step of the above, and received an acoustic model referred to in the case of speech recognition.

[0061]A recording medium which recorded a voice recognition program concerning this invention and in which computer reading is possible, A sound data acquisition function which acquires sound data which inputted an audio signal, performed speech recognition with reference to an acoustic model which asks for probability of an audio acoustical observed value series, makes realize a verification function which outputs a recognition result, and was updated, An acoustic model construction function to read updated sound data which the above-mentioned sound data acquisition function acquired in response to renewal instructions of an acoustic model, and to build an acoustic model, An acoustic model transmitting function which transmits an acoustic model which the above-mentioned acoustic model construction function built via a network, An acoustic model which the above-mentioned acoustic model transmitting function transmitted is received, and an acoustic model update function which the above-mentioned verification function updates by an acoustic model which received an acoustic model referred to in the case of speech recognition is realized.

[0062]A recording medium which recorded a voice recognition program concerning this invention and in which computer reading is possible, A language

data acquisition function which acquires language data which inputted an audio signal, performed speech recognition with reference to a language model which searches for appearing probability of a word sequence, makes realize a verification function which outputs a recognition result, and was updated, A language model construction function to read updated language data which the above-mentioned language data acquisition function acquired in response to renewal instructions of a language model, and to build a language model, A language model transmitting function which transmits a language model which the above-mentioned language model construction function built via a network, A language model which the above-mentioned language model transmitting function transmitted is received, and a language model update function which the above-mentioned verification function updates by a language model which received a language model referred to in the case of speech recognition is realized.

[0063]A recording medium which recorded a voice recognition program concerning this invention and in which computer reading is possible, A sound data acquisition function which acquires sound data which inputted an audio signal, performed speech recognition with reference to an acoustic model which asks for probability of an audio acoustical observed value series, makes realize a verification function which outputs a recognition result, and was updated, An updating acoustic model ID acquisition function which acquires ID which specifies the above-mentioned acoustic model in response to renewal instructions of an acoustic model, It corresponds to a specific condition directed by ID which the above-mentioned updating acoustic model ID acquisition function acquired, A sound data read-out function for the specification which reads updated sound data which the above-mentioned sound data acquisition function acquired, An acoustic model construction function for the specification which builds an acoustic model depending on the above-mentioned specific condition with reference to updated sound data which the above-mentioned sound data read-out function for specification read, An acoustic model transmitting function which transmits an acoustic model which the above-mentioned acoustic model construction function for specification built via a network, An acoustic model which the above-mentioned acoustic model transmitting function transmitted is received, and an acoustic model update function which the above-mentioned verification function updates by an acoustic model which received an acoustic model referred to in the case of speech recognition is realized.

[0064]A recording medium which recorded a voice recognition program concerning this invention and in which computer reading is possible, A language data acquisition function which acquires language data which inputted an audio signal, performed speech recognition with reference to a language model which searches for appearing probability of a word sequence, makes realize a verification function which outputs a recognition result, and was updated, An updating language model ID acquisition function which acquires ID which specifies the above-mentioned language model in response to renewal instructions of a language model, It corresponds to a specific condition directed by ID which the above-mentioned updating language model ID acquisition function acquired, A language data read-out function for the specification which reads updated language data which the above-mentioned language data acquisition function acquired, A language model construction function for the specification which builds a language model depending on the above-mentioned specific condition with reference to updated language data which the above-mentioned language data read-out function for specification read, A language model transmitting function which transmits a language model which the above-mentioned language model construction function for specification built via a network, A language model which the above-mentioned language model transmitting function transmitted is received, and a language model update function which the above-mentioned verification function updates by a language model which received a language model referred to in the case of speech recognition is realized.

[0065]A recording medium which recorded a voice recognition program concerning this invention and in which computer reading is possible, Input an audio signal and speech recognition is performed with reference to a language model which searches for appearing probability of a word sequence, and a user dictionary which registered a word, A language data acquisition function which acquires language data which makes realize a verification function which outputs a recognition result, and was updated, A user dictionary read-out function which reads the above-mentioned user dictionary in response to renewal instructions of a language model, A user dictionary dependence language model construction function to build a language model depending on a user dictionary which read updated language data which the above-mentioned language data acquisition function acquired, and the above-mentioned user dictionary read-out function read, A language model transmitting function which transmits a language model which the above-mentioned user dictionary dependence language model construction function built via a network, A language model which the above-mentioned language model transmitting function transmitted is received, and a language model update function which the above-mentioned verification function updates by a language model which received a language model referred to in the case of speech recognition is realized.

[0066]A recording medium which recorded a voice recognition program concerning this invention and in which computer reading is possible, A language data acquisition function which acquires language data which inputted an audio signal, performed speech recognition with reference to a language model which searches for appearing probability of a word sequence, makes realize a verification function which outputs a recognition result, and was updated, A user use text acquisition function which acquires a text which a user who performs speech recognition used in response to renewal instructions of a language model, A user use text dependence language model construction function to build a language model depending on a text which read updated language data which the above-mentioned language data acquisition function acquired, and the above-mentioned user use text acquisition function acquired, A language model transmitting function which transmits a language model which the above-mentioned user use text dependence language model construction function built via a network, A language model which the above-mentioned language model transmitting function transmitted is received, and a language model update function which the above-mentioned verification function updates by a language model which received a language model referred to in the case of speech recognition is realized.

[0067]A recording medium which recorded a voice recognition program concerning this invention and in which computer reading is possible, Input an audio signal and speech recognition is performed with reference to an acoustic model which asks for probability of an audio acoustical observed value series, An acoustic model ID acquisition function which acquires ID which realizes a verification function which outputs a recognition result and specifies the above-mentioned acoustic model, Read ID which the above-mentioned acoustic model ID acquisition function acquired, acquire voice data for adaptation from an inputted audio signal, and via a network, A voice acquisition function for adaptation which transmits voice data read ID and for acquired adaptation, Voice data for adaptation which the above-mentioned voice acquisition function for adaptation transmitted is used, An acoustic model adaptation function to carry out adaptation of the initial acoustic model before adaptation, and to match and store an adaptation finishing acoustic model in ID which the above-mentioned voice acquisition function for adaptation transmitted, In response to renewal instructions of an acoustic model, ID which the above-mentioned acoustic model ID acquisition function acquired via a network is received, An adaptation finishing acoustic model function preselection capability which chooses and reads an adaptation finishing acoustic model corresponding to ID which received out of an adaptation finishing acoustic model which the above-mentioned acoustic model adaptation function stored, An acoustic model transmitting function which transmits an adaptation finishing acoustic model which the above-mentioned adaptation finishing acoustic model function preselection capability read via a network, An adaptation finishing acoustic model which the above-mentioned acoustic model transmitting function transmitted is received, and an acoustic model update function which the above-mentioned verification function updates by an adaptation finishing acoustic model which received an acoustic model referred to in the case of speech recognition is realized.

[0068]

[Embodiment of the Invention]Hereafter, one gestalt of implementation of this invention is explained.

Embodiment 1. drawing 1 is a block diagram showing the composition of the voice recognition system by this embodiment of the invention 1. In a figure, the voice recognition equipment in which 10 performs speech recognition, the acoustic model managing server with which 20 is connected to the network, and 30 are language model managing servers connected to the network. Here, a network shows the general communication path which can transmit a digital signal by a cable or radio.

[0069]The audio signal which 100 inputs in the voice recognition equipment 10, the collation means to which 101 carries out speech recognition of the audio signal 100, The acoustic model in which the collation means 101 refers to 102 in the case of speech recognition, the language model in which the collation means 101 refers to 103 in the case of speech recognition, The recognition result to which the collation means 101 outputs 104, the acoustic model update means which updates the acoustic model 102 by the acoustic model to which 111 was transmitted via the network, and 118 are language

model update means which update the language model 103 by the language model transmitted via the network.

[0070]The renewal instructions of an acoustic model to which 105 is given from the exterior in the acoustic model managing server 20, The sound data acquisition means which acquires the sound data in which 106 was updated, the updated sound data in which the sound data acquisition means 106 acquired 107, 108 reads the updated sound data 107 in response to the renewal instructions 105 of an acoustic model, The acoustic model construction means which performs parameter estimation using a statistical method and builds an acoustic model, The acoustic model storing means in which 109 stores the built acoustic model, and 110 are acoustic model transmitting means which transmit the acoustic model stored in the acoustic model storing means 109 to the voice recognition equipment 10 via a network.

[0071]The renewal instructions of a language model to which 112 is given from the exterior in the language model managing server 30, The language data acquisition means which acquires the language data in which 113 was updated, the updated language data in which the language data acquisition means 113 acquired 114, 115 reads the updated language data 114 in response to the renewal instructions 112 of a language model, The language model construction means which performs parameter estimation using a statistical method and builds a language model, The language model storing means in which 116 stores the built language model, and 117 are language model transmitting means which transmit the language model stored in the language model storing means 116 to the voice recognition equipment 10 via a network.

[0072]The characteristic portion of this different invention from conventional technology, In the acoustic model managing server 20 and the language model managing server 30, by the sound data acquisition means 106 and the language data acquisition means 113. The updated sound data 107 and the updated language data 114 are acquired, By the acoustic model construction means 108 and the language model construction means 115, the newest acoustic model, In [ build the newest language model, transmit the built newest acoustic model and the newest language model to the voice recognition equipment 10 via a network, and ] the voice recognition equipment 10, The acoustic model update means 111 and the language model update means 118 are updating the acoustic model 102 which the collation means 101 refers to, and the language model 103 by the newest acoustic model and the newest language model.

[0073]Next, operation is explained. In the acoustic model managing server 20, it operates asynchronously, and the sound data acquisition means 106 downloads automatically or semi-automatically the renewal instructions 105 of an acoustic model, a synchronization, or the sound data always updated or distributed, and stores it in the updated sound data 107. These sound data to acquire is updated, for example on the Internet, or, It is the voice data distributed by multimedia broadcast or voice data, and a corresponding opening sentence text, and it is searched on the Internet by a retrieving tool, or determines and downloads using the race card of multimedia broadcast.

[0074]The updated sound data 107 is accumulation of the sound data for acoustic model study acquired by the sound data acquisition means 106, and consists of voice data, voice data, and a corresponding opening sentence text in the above-mentioned example.

[0075]A time interval with the acoustic model construction means 108 constant, for example, the time interval with which speech recognition processing was carried out, Or the renewal instructions 105 of an acoustic model given to suitable timing, such as a user's directions given from an input device, are received, With reference to the updated sound data 107, by performing parameter estimation of an acoustic model using a statistical method, For example, by using a vector quantization algorithm only from voice data, Or by using a BAUMU Welch algorithm from voice data and a corresponding opening sentence text, an acoustic model is built and it stores in the acoustic model storing means 109 so that learned data may be expressed well.

[0076]The acoustic model storing means 109 memorizes the acoustic model built by the acoustic model construction means 108, and outputs an acoustic model according to a read request. The acoustic model transmitting means 110 reads an acoustic model from the acoustic model storing means 109, and transmits to the acoustic model update means 111 of the voice recognition equipment 10 via a network.

[0077]In the voice recognition equipment 10, the acoustic model update means 111 updates the acoustic model 102 referred to in the case of collation of the collation means 101 by the acoustic model received from the acoustic model transmitting means 110 of the acoustic model managing server 20 via the network.

[0078]In the language model managing server 30, the language data acquisition means 113, It operates asynchronously and the renewal instructions 112 of a language model, a synchronization, or the language data always updated or distributed is downloaded, and the new language data used for language model construction is collected, and is stored in the updated language data 114. These language data to acquire is a text, texts, such as a chat, mail, and a manual, etc. which can be searched from on the newspaper and mail magazine which are distributed periodically, for example, or the Internet.

[0079]The updated language data 114 is accumulation of the language model study declinable word word data acquired by the language data acquisition means 113, and is text data, the keyword information about the contents of a text acquired simultaneously, etc.

[0080]A time interval with the language model construction means 115 constant, for example, the time interval with which speech recognition processing was carried out, Or the renewal instructions 112 of a language model given to suitable timing, such as a user's directions given from an input device, are received, By reading text data with reference to the updated language data 114, and performing parameter estimation of a language model using a statistical method, For example, by calculating n gram statistics value from the text data divided into the word, a language model is built so that the data for study may be expressed well, and it stores in the language model storing means 116.

[0081]The language model storing means 116 memorizes the language model built by the language model construction means 115, and outputs a language model according to a read request. The language model transmitting means 117 reads a language model from the language model storing means 116, and transmits to the language model update means 118 of the voice recognition equipment 10 via a network.

[0082]In the voice recognition equipment 10, the language model update means 118 updates the language model 103 referred to in the case of collation of the collation means 101 by the language model received from the language model transmitting means 117 of the language model managing server 30 via the network.

[0083]Drawing 2 is a flow chart which shows the update process of the acoustic model 102 by this embodiment of the invention 1. In step ST201, an acoustic model update timing determination means (not shown), For example, suitable update timing is judged from the surveillance of the using state of the time interval from the last modification time of a user's directions or an acoustic model, and a network, etc., and the renewal instructions 105 of an acoustic model are transmitted to the acoustic model construction means 108 of the acoustic model managing server 20. If the renewal instructions 105 of an acoustic model are received, the acoustic model construction means 108 progresses to step ST202, and if it has not received the renewal instructions 105 of an acoustic model, it will end processing.

[0084]In step ST202, the acoustic model construction means 108 reads the updated sound data 107 which is used for study. In step ST203, by performing parameter estimation of an acoustic model using a statistical method, the acoustic model construction means 108 builds an acoustic model, and stores the built acoustic model in the acoustic model storing means 109 from the updated sound data 107.

[0085]In step ST204, the acoustic model transmitting means 110 reads an acoustic model from the acoustic model storing means 109, and transmits an acoustic model to the acoustic model update means 111 of the voice recognition equipment 10 via a network. In step ST205, the acoustic model update means 111 updates the acoustic model 102 which the collation means 101 refers to by the received acoustic model.

[0086]When requiring renewal of the acoustic model 102 by the renewal instructions 105 of an acoustic model, The version of the acoustic model 102 which the voice recognition equipment 10 uses simultaneously is transmitted, If the acoustic model transmitting means 110 transmits only difference information with not the whole acoustic model stored in the acoustic model storing means 109 but the acoustic model 102 which the voice recognition equipment 10 used till then, send data can be reduced and network load can be reduced.

[0087]The acoustic model construction means 108 builds the acoustic model updated beforehand, and even when a gestalt which transmits an acoustic model to the voice recognition equipment 10 is taken according to the demand of the renewal instructions 105 of an acoustic model, it can operate similarly.

[0088]Even if it is a case where the voice recognition equipment 10 has a user dictionary, it can process similarly.



[0089]In this embodiment, although it explains for speech recognition, if it is aimed at the pattern recognition which consists of a stochastic model showing the relation between a pattern and a symbol, and a stochastic model showing the appearance of the symbol, it is applicable similarly.

[0090]As long as the storing form of the updated sound data 107 is a form available at the time of acoustic model construction, it may calculate signal processing and frequency distribution beforehand.

[0091]Drawing 3 is a flow chart which shows the update process of the language model 103 by this embodiment of the invention 1. In step ST301, a language model update timing determination means (not shown), For example, suitable update timing is judged from the surveillance of the using state of the time interval from the last modification time of a user's directions or a language model, and a network, etc., and the renewal instructions 112 of a language model are transmitted to the language model construction means 115 of the language model managing server 30. If the renewal instructions 112 of a language model are received, and the language model construction means 115 progresses to step ST302 and has not received the renewal instructions 112 of a language model, it ends processing.

[0092]In step ST302, the language model construction means 115 reads the updated language data 114 which is used for study. In step ST303, by performing parameter estimation of a language model using a statistical method, the language model construction means 115 builds a language model, and stores the built language model in the language model storing means 116 from the updated language data 114.

[0093]In step ST304, the language model transmitting means 117 reads a language model from the language model storing means 116, and transmits to the language model update means 118 of the voice recognition equipment 10 via a network. In step ST305, the language model update means 118 updates the language model 103 which the collation means 101 refers to by the received language model.

[0094]By transmitting the version of the language model 103 which the voice recognition equipment 10 uses simultaneously by the renewal instructions 112 of a language model when requiring renewal of a language model, If only difference information with not the whole language model stored in the language model storing means 116 but the language model 103 which the voice recognition equipment 10 used till then is transmitted, send data can be reduced and network load can be reduced.

[0095]The language model construction means 115 builds the language model updated beforehand, and even when a gestalt which transmits a language model is taken according to the demand of the renewal instructions 112 of a language model, it can operate similarly.

[0096]Even if it is a case where the voice recognition equipment 10 has a user dictionary, it can operate similarly.

[0097]In this embodiment, although it explains for speech recognition, if it is aimed at the pattern recognition which consists of a stochastic model showing the relation between a pattern and a symbol, and a stochastic model showing the appearance of the symbol, it is applicable similarly.

[0098]The storing form of the updated language data 114, if it is a form available at the time of language model construction, it can divide into the word beforehand or can use for language model construction — as — a word, a word chain, the combination of the word which appears simultaneously, etc. — frequency — or probability calculation may be carried out.

[0099]Although drawing 1 of this Embodiment 1 shows the case where an acoustic model is built according to the renewal instructions 105 of an acoustic model, When the sound data acquisition means 106 acquires sound data, the acoustic model construction means 108 updates an acoustic model, and it stores in an acoustic model storing means 109, and may be made for the acoustic model transmitting means 110 to read an acoustic model in response to the renewal instructions 105 of an acoustic model. The same may be said of renewal of a language model.

[0100]Although drawing 1 showed the case where it had the acoustic model update means 111 and the language model update means 118, you may be a case where it has either the acoustic model update means 111 or the language model update means 118.

[0101]It is also recordable on a recording medium by making the voice recognition system in Embodiment 1 into a voice recognition program. In this case, the sound data acquisition function which performs the same processing as the sound data acquisition means 106 in the acoustic model managing server 20, The acoustic model construction function to perform the same processing as the acoustic model construction means 108, In the software which comprises an acoustic model storing function to perform the same processing as the acoustic model storing means 109, and an acoustic model transmitting function which performs the same processing as the acoustic model transmitting means 110, and the language model managing server 30, The language data acquisition function which performs the same processing as the language data acquisition means 113, The language model construction function to perform the same processing as the language model construction means 115, In the software which comprises a language model storing function to perform the same processing as the language model storing means 116, and a language model transmitting function which performs the same processing as the language model transmitting means 117, and the voice recognition equipment 10, It becomes a voice recognition program by the software which comprises an acoustic model update function which performs the same processing as the acoustic model update means 111, a language model update function which performs the same processing as the language model update means 118, and a verification function which performs the same processing as the collation means 101.

[0102]The voice recognition program recorded on a recording medium The software of the voice recognition equipment 10, The software of the acoustic model managing server 20 may be recorded on a separate recording medium, may be recorded on one recording medium, and it may transmit to the acoustic model managing server 20 or the voice recognition equipment 10, respectively from the voice recognition equipment 10 or the acoustic model managing server 20. This is the same even when a language model is targeted.

[0103]With as mentioned above, the acoustic model managing server 20 or the language model managing server 30 which was connected to the network according to this Embodiment 1. The updated sound data 107 or the updated language data 114 is acquired, By building the acoustic model or language model in the newest state, and updating users' acoustic model 102 or language model 103 of the voice recognition equipment 10 via a network. The effect that the recognition precision of speech recognition can be raised is acquired without applying a big burden to a user.

[0104]Embodiment 2. drawing 4 is a block diagram showing the composition of the voice recognition system by this embodiment of the invention 2. The updating language model ID acquiring means which acquires ID which specifies the language model 103 of the voice recognition equipment 10 which 401 updates in the language model managing server 30 of drawing 4, and 402, Corresponding to the specific condition directed by ID which updating language model ID acquiring means 401 acquired, the language data reading means for the specification which reads the updated language data 114, and 403, It is a language model construction means for the specification which builds the language model corresponding to a specific condition with reference to the updated language data 114 which was read by the language data reading means 402 for specification. About each means and each model which were already explained, the same numerals are attached and explanation is omitted.

[0105]A portion characteristic of this different embodiment from conventional technology is providing the language model which was provided with updating language model ID acquiring means 401, the language data reading means 402 for specification, and the language model construction means 403 for specification, and was specified by language model ID via the network. Here, the language model for specification is a language model learned so that higher performance might be obtained by making a language model specialize in a specific user, a group, application application, etc.

[0106]Next, operation is explained. In the language managing server 30, updating language model ID acquiring means 401 acquires ID used in order to choose the specific language model 103 from two or more updated language models 103. This ID is a user's user ID, ID showing the task which is the target of the audio signal 100, etc., and can provide the language model 103 for the specification to update in a meaning, for example.

[0107]The language data reading means 402 for specification receives updating language model ID which updating language model ID acquiring means 401 acquired, and reads the updated language data 114 in a sentence or the independent unit of a text, For example, it judges from the keyword about the contents of a text given to language data, or the keyword contained in language data, the flag which identifies whether it is an object specified by language model ID is given, and it sends to the language model construction means 403 for specification.

[0108]The language model construction means 403 for specification builds the language model learned from the updated language data 114 so that recognition precision might become high about a specific object, and stores it in the language model storing means 116.

[0109]For this reason, first study declinable word word data from the keyword etc. which are contained considering a sentence or two or more sentences as a unit. Judge the language model 103 specified and according to this, for example "Examination of the prior task adaptation for interactive

voice recognition", it inquires in Akinori Ito, Masanori Koda, the Institute of Electronics, Information and Communication Engineers technical research report (SP96-81), and 1996 (it is hereafter considered as the literature 4) -- as -- relation -- the language model corresponding to a specific condition can be built by giving big dignity by deep text data.

[0110]For example, in order to build the language model for the specification which specialized in the topic about a sport, With reference to the flag obtained from the language data reading means 402 for specification at the time of language model study, if it is text data of the topic about a sport, actual frequency will be carried out alpha double and will be counted, if it is the other report, both will be added by frequency as they are, and a stochastic model will be presumed. Here, alpha determines that the entropy of the language model to the data which is not used for study among the texts for the specification made into the object of speech recognition serves as the minimum.

[0111]Drawing 5 is a flow chart which shows the update process of the language model 103 by this embodiment of the invention 2. In step ST501, a language model update timing determination means (not shown), For example, suitable update timing is judged from the surveillance of the using state of the time interval from a user's directions or the last modification time, and a network, etc., and the renewal instructions 112 of a language model are transmitted to updating language model ID acquiring means 401 of the language model managing server 30. If the renewal instructions 112 of a language model are received, and updating language model ID acquiring means 401 progresses to step ST502 and has not received the renewal instructions 112 of a language model, it ends processing.

[0112]In step ST502, if it is update timing, by means to specify the user group who is using it, a means to specify a task, etc., updating language model ID acquiring means 401 will acquire ID of the language model 103 of an updating place, and will send it to the language data reading means 402 for specification. In step ST503, the language data reading means 402 for specification, According to language model ID for specification, the updated language data 114 is read in a sentence or the independent unit of a text, the flag which distinguishes whether it is an object specified by language model ID is given, and the updated language data 114 is read.

[0113]In step ST504, according to learning algorithm, the language model construction means 403 for specification presumes a language model from the updated language data 114, and stores the presumed language model in the language model storing means 116. In step ST505, the language model transmitting means 117 transmits the read language model to the language model update means 118 of the voice recognition equipment 10 via a network. In step ST506, the language model update means 118 updates the language model 103 which the collation means 101 refers to by the received language model.

[0114]If the suiting language model can be outputted to language model ID, it is not necessary to build the language model for specification beforehand. For example, only study declinable word word data is created depending on language model ID, and a language model may be built according to a demand.

[0115]Although language model construction for specification according to the literature 4 was made into the example in this embodiment, if it is the method of choosing from two or more language models 103, it is applicable similarly using ID which can determine a language model.

[0116]By transmitting the version of the language model 103 which the voice recognition equipment 10 uses simultaneously in the case of a language model update request by the renewal instructions 112 of a language model, Only the difference information from not the whole language model for the built specification but the present language model 103 can be transmitted, and network load can be reduced.

[0117]Even if it is a case where voice recognition equipment has the user dictionary 601, it can process similarly.

[0118]Although this explanation explained for speech recognition, if aimed at the pattern recognition which consists of a stochastic model showing the relation between a pattern and a symbol, and a stochastic model showing the appearance of the symbol, it is applicable similarly.

[0119]Although the language model for specification was built and the language model 103 is updated by the language model for the built specification in this embodiment, the acoustic model for specification can be built and the acoustic model 102 can also be updated by the acoustic model for the built specification. In drawing 4, instead of the language model managing server 30, instead of an acoustic model managing server and the language data acquisition means 113 In that case, a sound data acquisition means, Instead of the sound data and updating language model ID acquiring means 401 which were updated instead of the updated language data 114, an updating acoustic model ID acquiring means, Instead of the language data reading means 402 for specification, the acoustic model reading means for specification, Instead of the language model construction means 403 for specification, the acoustic model construction means for specification, In [ have an acoustic model transmitting means instead of the language model storing means 116 instead of an acoustic model storing means and the language model transmitting means 117, and ] the voice recognition equipment 10, It has an acoustic model update means instead of the language model update means 118, and an acoustic model update means should just update the acoustic model 102.

[0120]It is also recordable on a recording medium by making the voice recognition system in Embodiment 2 into a voice recognition program. In this case, the language data acquisition function which performs the same processing as the language data acquisition means 113 in the language model managing server 30, The updating language model ID acquisition function which performs the same processing as updating language model ID acquiring means 401, The language data read-out function for the specification which performs the same processing as the language data reading means 402 for specification, The language model construction function for the specification which performs the same processing as the language model construction means 403 for specification, In the software which comprises a language model storing function to perform the same processing as the language model storing means 116, and a language model transmitting function which performs the same processing as the language model transmitting means 117, and the voice recognition equipment 10, It becomes a voice recognition program by the software which comprises a language model update function which performs the same processing as the language model update means 118, and a verification function which performs the same processing as the collation means 101. This is the same even when an acoustic model is targeted.

[0121]The voice recognition program recorded on a recording medium The software of the voice recognition equipment 10, The software of the language model managing server 30 may be recorded on a separate recording medium, may be recorded on one recording medium, and it may transmit to the language model managing server 30 or the voice recognition equipment 10, respectively from the voice recognition equipment 10 or the language model managing server 30. This is the same even when an acoustic model is targeted.

[0122]With as mentioned above, the language model managing server 30 which was connected to the network according to this Embodiment 2. By acquiring the updated language data 114 and moreover acquiring ID of the language model 103 a user's voice recognition equipment 10. By being in the newest state, building the language model for the specification corresponding to a user moreover, and updating the language model 103 of a user's voice recognition equipment 10 via a network. The effect that the recognition precision of speech recognition can be raised is acquired without applying a big burden to a user.

[0123]By updating the language model customized for specification via a network according to this Embodiment 2, Even when a user uses several different collation means 101, the suitable language model 103 for the utilization time of all the collation means 101 is used, and the effect that high recognition precision can be acquired is acquired.

[0124]Embodiment 3. drawing 6 is a block diagram showing the composition of the voice recognition system by this embodiment of the invention 3. In [ in the voice recognition equipment 10 of drawing 6 601 is the user dictionary which registered the word referred to in the case of collation of the collation means 101, and ] the language model managing server 30, The user dictionary reading means which reads the user dictionary 601 with which the collation means 101 refers to 602 in response to the renewal instructions 112 of a language model via a network, and 603, It is a user dictionary dependence language model construction means which builds the language model depending on the user dictionary 601 with reference to the updated language data 114 and the user dictionary 601 which the user dictionary reading means 602 read. About each means and each model which were already explained, the same numerals are attached and explanation is omitted.

[0125]A portion characteristic of this different embodiment from conventional technology is having had the user dictionary reading means 602 and the user dictionary dependence language model construction means 603.



[0126]Next, operation is explained. The user dictionary reading means 602 of the language model managing server 30 reads the user dictionary 601 which the collation means 101 of the voice recognition equipment 10 refers to in response to the renewal instructions 112 of a language model via a network. The user dictionary dependence language model construction means 603 builds the language model which is updated by the newest state and customized by the user using the word registered into the user dictionary 601, and the updated language data 114.

[0127]Construction of the language model depending on the user dictionary 601, For example, the text in which the word which exists in the user dictionary 601 is contained out of the updated language data 114 is extracted, and it considers that this is a text for specification, and is carried out by enforcing the method of the literature 4 statement referred to by Embodiment 2. It is expectable that it is not registered by the original language model among the words of user dictionary 601 statement, but become possible to give a statistics value appropriate to the word which appeared in the updated text with the word to which the suitable statistics value was not given, and recognition precision improves by this.

[0128]Drawing 7 is a flow chart which shows the update process of the language model 103 by this embodiment of the invention 3. In step ST701, a language model update timing determination means (not shown), For example, suitable update timing is judged from the surveillance of the using state of the time interval from a user's directions or the last modification time, and a network, etc., and the renewal instructions 112 of a language model are transmitted to the user dictionary reading means 602 of the language model managing server 30. If the renewal instructions 112 of a language model are received, and the user dictionary reading means 602 progresses to step ST702 and has not received the renewal instructions 112 of a language model, it ends processing.

[0129]In step ST702, the user dictionary reading means 602 reads the user dictionary 601 which the collation means 101 refers to via a network. In step ST703, the user dictionary dependence language model construction means 603 reads the language data 114 updated further. In step ST704, the user dictionary dependence language model construction means 603 builds the language model which was dependent on the user dictionary 601 from the user dictionary 601 and the updated language data 114, and stores it in the language model storing means 116.

[0130]In step ST705, the language model transmitting means 117 transmits the language model depending on the user dictionary 601 read from the language model storing means 116 to the language model update means 118 of the voice recognition equipment 10 via a network. In step ST706, the language model update means 118 updates the language model 103 which the collation means 101 refers to by the received language model.

[0131]In this embodiment, although it explains for speech recognition, if it is aimed at the pattern recognition which consists of a stochastic model showing the relation between a pattern and a symbol, and a stochastic model showing the appearance of the symbol, it is applicable similarly.

[0132]It is also recordable on a recording medium by making the voice recognition system in Embodiment 3 into a voice recognition program. In this case, the language data acquisition function which performs the same processing as the language data acquisition means 113 in the language model managing server 30, The user dictionary read-out function to perform the same processing as the user dictionary reading means 602, The user dictionary dependence language model construction function to perform the same processing as the user dictionary dependence language model construction means 603, In the software which comprises a language model storing function to perform the same processing as the language model storing means 116, and a language model transmitting function which performs the same processing as the language model transmitting means 117, and the voice recognition equipment 10, It becomes a voice recognition program by the software which comprises a language model update function which performs the same processing as the language model update means 118, and a verification function which performs the same processing as the collation means 101.

[0133]The voice recognition program recorded on a recording medium The software of the voice recognition equipment 10, The software of the language model managing server 30 may be recorded on a separate recording medium, may be recorded on one recording medium, and it may transmit to the language model managing server 30 or the voice recognition equipment 10, respectively from the voice recognition equipment 10 or the language model managing server 30.

[0134]With as mentioned above, the language model managing server 30 which was connected to the network according to this Embodiment 3. Acquire the updated language data 114 and, moreover, by reading the user dictionary 601 of a user's voice recognition equipment 10. By being in the newest state, building the language model made to reflect in details more about the word moreover registered into the user dictionary 601, and updating the language model 103 of a user's voice recognition equipment 10 via a network. The effect that the recognition precision of speech recognition can be raised is acquired without applying a big burden to a user, even when a user dictionary becomes large.

[0135]Embodiment 4. drawing 8 is a block diagram showing the composition of the voice recognition system by this embodiment of the invention 4. In the language model managing server 30 of drawing 8, 801, The user use text acquisition means which acquires the text which the user used in response to the renewal instructions 112 of a language model, The user use text storing means in which 802 stores the text which the user use text acquisition means 801 acquired, and 803, It is a user use text dependence language model construction means which builds the language model depending on a text with reference to the updated language data 114 and the text stored in the user use text storing means 802. About each means and each model which were already explained, the same numerals are attached and explanation is omitted.

[0136]A portion characteristic of this different embodiment from conventional technology, It has the user use text acquisition means 801, the user use text storing means 802, and the user use text dependence language model construction means 803, It is building a language model according to the text which the user used, and the text which the user used with reference to the language data 114 updated by the newest state.

[0137]Next, operation is explained. The user use text acquisition means 801 of the language model managing server 30 reads the text file which the user referred to or described by scanning the file and directory which the user specified beforehand, in response to the fact that the renewal instructions 112 of a language model. The user use text storing means 802 stores the text collected by the user use text acquisition means 801.

[0138]With reference to a user use text and the updated language data 114, the user use text dependence language model construction means 803 builds a language model so that recognition precision may become high. In construction of the language model using a user use text, the language model of user use text dependence is built by considering that a user use text is a text for specification for example, and enforcing the method of the literature 4 statement referred to by Embodiment 2. thus, since the user is making the character of the text which carried out reference or existing appearance reflect, the built language model can obtain a higher-precision recognition result including the linguistic character in which the probability which a user utters is high.

[0139]Drawing 9 is a flow chart which shows the update process of the language model 103 by this embodiment of the invention 4. In step ST901, a language model update timing determination means (not shown), Suitable timing is judged from the monitor of the using state of the time interval from a user's directions or the last modification time, and a network, etc., and the renewal instructions 112 of a language model are transmitted to the user use text acquisition means 801 of the language model managing server 30. If the renewal instructions 112 of a language model are received, and the user use text acquisition means 801 progresses to step ST902 and has not received the renewal instructions 112 of a language model, it ends processing.

[0140]In step ST902, the user use text acquisition means 801 reads a user use text, and stores it in the user use text storing means 802. In step ST903, the user use text dependence language model construction means 803 reads a user use text and the updated language data 114. In step ST904, the user use text dependence language model construction means 803 builds a user use text dependence language model from a user use text and the updated language data 114, and stores it in the language model storing means 116.

[0141]In step ST905, the language model transmitting means 117 transmits the user use text dependence language model read from the language model storing means 116 to the language model update means 118 of the voice recognition equipment 10 via a network. In step ST906, the language model update means 118 updates the language model 103 which the collation means 101 refers to by the received language model.

[0142]Although it presupposed that acquisition of a user use text is based on search of a specific directory or a file in this embodiment, The user use text acquisition means 801 will pick out a text from neither a file nor a directory, if a text is collectable, The text etc. which the user perused by user inputs, such as speech recognition, a keyboard, a pen, and OCR, or a browser may be used.

[0143]Although the user use text storing means 802 presupposed that the text collected by the user use text acquisition means 801 is stored, It is the

same, even if a suitable means divides a text into a word or it stores in accordance with the standard of the user use text dependence language model construction means 803 as frequency about the word referred to at the time of modern construction, a word chain, the combination of the word which appears simultaneously, etc.

[0144]In this embodiment, although explained for speech recognition, if aimed at the pattern recognition which consists of a stochastic model showing the relation between a pattern and a symbol, and a stochastic model showing the appearance of the symbol, it is applicable similarly.

[0145]The voice recognition system in Embodiment 4 is also recordable on a recording medium as a voice recognition program. In this case, the language data acquisition function which performs the same processing as a language data acquisition means in the language model managing server 30. The user use text acquisition function which performs the same processing as the user use text acquisition means 801, The user use text storing function to perform the same processing as the user use text storing means 802, The user use text dependence language model construction function to perform the same processing as the user use text dependence language model construction means 803, In language model storing which performs the same processing as the language model storing means 116, the software which comprises a language model transmitting function which performs the same processing as the language model transmitting means 117, and the voice recognition equipment 10, It becomes a voice recognition program by the software which comprises a language model update function which performs the same processing as the language model update means 118, and a verification function which performs the same processing as the collation means 101.

[0146]The voice recognition program recorded on a recording medium The software of the voice recognition equipment 10, The software of the language model managing server 30 may be recorded on a separate recording medium, may be recorded on one recording medium, and it may transmit to the language model managing server 30 or the voice recognition equipment 10, respectively from the voice recognition equipment 10 or the language model managing server 30.

[0147]With as mentioned above, the language model managing server 30 which was connected to the network according to this Embodiment 4. By acquiring the updated language data 114 and moreover acquiring the text which the user used. By being in the newest state, building the language model depending on the text which the user moreover used, and updating the language model 103 of a user's voice recognition equipment 10 via a network. The effect that the recognition precision of speech recognition can be raised is acquired without applying a big burden to a user.

[0148]Embodiment 5. drawing 10 is a block diagram showing the composition of the voice recognition system by this embodiment of the invention 5. The acoustic model ID acquiring means which acquires ID which identifies the acoustic model which refers to 1001 in the case of collation of the collation means 101 in the voice recognition equipment 10 of drawing 10, and 1002, It is a voice acquisition means for adaptation which reads ID which acoustic model ID acquiring means 1001 acquired, acquires the voice data for adaptation from the inputted audio signal, and transmits read ID and the acquired voice data for adaptation to the acoustic model managing server 20 via a network.

[0149]1003 in the acoustic model managing server 20 of drawing 10 the initial acoustic model before adaptation and 1004, The voice data for adaptation transmitted from the voice acquisition means 1002 for adaptation of the voice recognition equipment 10 is used, Carry out adaptation of the initial acoustic model 1003 before adaptation, and an adaptation finishing acoustic model, An acoustic model adaptation means to match with ID transmitted from the voice acquisition means 1002 for adaptation, and to store in the adaptation finishing acoustic model storing means 1005, and 1006, In response to the renewal instructions 105 of an acoustic model, ID which acoustic model ID acquiring means 1001 of the voice recognition equipment 10 acquired via the network is received, It is an adaptation finishing acoustic model selecting means which chooses and reads the adaptation finishing acoustic model corresponding to ID which received from the adaptation finishing acoustic model storing means 1005. About each means and each model which were already explained, the same numerals are attached and explanation is omitted.

[0150]A portion characteristic of this different embodiment from conventional technology, In [ in the voice recognition equipment 10, have acoustic model ID acquiring means 1001 and the voice acquisition means 1002 for adaptation, and ] the acoustic model managing server 20, It is having had the initial acoustic model 1003, the acoustic model adaptation means 1004, the adaptation finishing acoustic model storing means 1005, and the adaptation finishing acoustic model selecting means 1006.

[0151]In [ in this embodiment, acquire the voice data acoustic model ID used as an adaptation object, and for adaptation in the voice recognition equipment 10, and ] the acoustic model managing server 20, In [ build the acoustic model which performed adaptation depending on acoustic model ID, transmit to the voice recognition equipment 10, and ] the voice recognition equipment 10, Since the user can refer to the acoustic model which carried out adaptation when using the arbitrary collation means 101 by updating the acoustic model 102, higher recognition precision can be acquired.

[0152]Next, operation is explained. In the voice recognition equipment 10, acoustic model ID acquiring means 1001 is the user ID of the user who determines the acoustic model used as an adaptation object, and uses the voice recognition equipment 10. Acoustic model ID which the voice acquisition means 1002 for adaptation acquires the voice data for adaptation based on the audio signal 100 as an object for adaptation beforehand before use of the collation means 101, and is read from acoustic model ID acquiring means 1001, The acquired bottom transmits the voice data for adaptation to an acoustic model adaptation means 1004 of the acoustic model managing server 20 by which it is connected via a network.

[0153]In the acoustic model managing server 20, the initial acoustic model 1003 is an acoustic model before performing adaptation. The voice data and the initial acoustic model 1003 for adaptation which were received via the network are used for the acoustic model adaptation means 1004, The acoustic model which carried out adaptation is built and acoustic model ID received via an acoustic model and a network is stored in the adaptation finishing acoustic model storing means 1005. [ finishing / adaptation ] A maximum-a-posteriori-probability-estimation method is used for the adaptation of an acoustic model, for example.

[0154]The adaptation finishing acoustic model storing means 1005 stores the acoustic model by which adaptation was carried out to acoustic model ID by the acoustic model adaptation means 1004, and outputs an acoustic model with appointed acoustic model ID according to a demand of an adaptation finishing acoustic model selecting means. In response to the renewal instructions 105 of an acoustic model, the adaptation finishing acoustic model selecting means 1006 acquires acoustic model ID from acoustic model ID acquiring means 1001, and chooses and reads a corresponding adaptation finishing acoustic model from the adaptation finishing acoustic model storing means 1005.

[0155]Drawing 11 is a flow chart which shows the update process of the acoustic model 102 by this embodiment of the invention 5. Processing is divided into the adaptation stage of the acoustic models from step ST1101 to ST1107, and the updating stage of the acoustic models from step ST1108 to ST1112. In an adaptation stage, the acoustic model depending on acoustic model ID is built using the inputted voice data for adaptation.

[0156]In step ST1101, acoustic model ID acquiring means 1001 of the voice recognition equipment 10 identifies the acoustic model used as an adaptation object, for example, acquires identification information, such as a user name. In step ST1102, the voice acquisition means 1002 for adaptation reads acoustic model ID, and acquires the voice data for adaptation inputted from the audio signal 100. In step ST1103, the voice acquisition means 1002 for adaptation transmits the adaptation demand of an acoustic model to the acoustic model managing server 20 via a network, and transmits the voice data for acoustic model ID and adaptation to the acoustic model adaptation means 1004 simultaneously.

[0157]In step ST1104, the acoustic model adaptation means 1004 receives the adaptation demand of an acoustic model via a network, and reads the voice data for acoustic model ID and adaptation. In step ST1105, the acoustic model adaptation means 1004 reads the initial acoustic model 1003. In step ST1106, the acoustic model adaptation means 1004 carries out adaptation of the initial acoustic model 1003 using the voice data for adaptation received via the network. In step ST1107, the acoustic model adaptation means 1004 is stored in the adaptation finishing acoustic model storing means 1005 so that the acoustic model by which adaptation was carried out can be distinguished by acoustic model ID.

[0158]In step ST1108, an acoustic model update timing determination means (not shown). Suitable update timing is judged from the surveillance of the using state of the time interval from a user's directions or the last modification time, and a network, etc., and the renewal instructions 105 of an acoustic model are transmitted to the adaptation finishing acoustic model selecting means 1006. If the renewal instructions 105 of an acoustic model are received, the adaptation finishing acoustic model selecting means 1006 progresses to step ST1109, and if it has not received the renewal instructions 105 of an acoustic model, it will end processing. In step ST1109, the adaptation finishing acoustic model selecting means 1006 reads

acoustic model ID for adaptation from acoustic model ID acquiring means 1001 of the voice recognition equipment 10 via a network.

[0159]In step ST1110, the adaptation finishing acoustic model selecting means 1006 chooses and reads the acoustic model specified by acoustic model ID from the adaptation finishing acoustic model storing means 1005. In step ST1111, the acoustic model transmitting means 110 transmits the read adaptation finishing acoustic model to the acoustic model update means 111 of the voice recognition equipment 10 via a network. In step ST1112, the acoustic model update means 111 updates the acoustic model 102 which the collation means 101 refers to by the received adaptation finishing acoustic model.

[0160]In the voice recognition equipment 10, as long as acoustic model ID acquiring means 1001 determines the acoustic model used as an adaptation object, it may determine a circuit transfer characteristic, the background-noise characteristic, the reverberant sound characteristic, etc.

[0161]In the voice recognition equipment 10, although the voice acquisition means 1002 for adaptation acquires a user's voice data as an object for adaptation beforehand before use of the collation means 101, It is also possible to perform adaptation of the acoustic model referred to in the case of the next collation in the collation means 101 acquiring and carrying out voice data at the time of collation, and inputting into the voice acquisition means 1002 for adaptation.

[0162]In the voice recognition equipment 10, the procedure of acquisition of the data of the sound for adaptation and acquisition of acoustic model ID may be reverse.

[0163]In the voice recognition equipment 10, although the storing gestalt of the voice data of the voice acquisition means 1002 for adaptation was made into the voice waveform, The time series of the sound feature vector which carried out signal processing of the voice data, the code book code sequence acquired with reference to a sound feature vector and a vector quantization code book, As long as the frequency distribution acquired by carrying out the statistical work of them, the probability distributions acquired from frequency distribution, etc. are the gestalten which can be used for study of an acoustic model, what kind of thing may be used.

[0164]A means to store the voice data for adaptation which the acoustic model adaptation means 1004 received in the acoustic model managing server 20 is added, By updating the early model 1003 with the stored voice data for much adaptation, the accuracy of study becomes high and the adaptation finishing acoustic model which performs highly precise recognition can be built.

[0165]Although speech recognition was targeted in this embodiment, if aimed at the pattern recognition which consists of a stochastic model showing the relation between a pattern and a symbol, and a stochastic model showing the appearance of the symbol, it is applicable similarly.

[0166]It is also recordable on a recording medium by making the voice recognition system in Embodiment 5 into a voice recognition program. In this case, the acoustic model ID acquisition function which performs the same processing as acoustic model ID acquiring means 1001 in the voice recognition equipment 10, The voice acquisition function for adaptation which performs the same processing as the voice acquisition means 1002 for adaptation, In the software which comprises an acoustic model update function which performs the same processing as the acoustic model update means 111, and a verification function which performs the same processing as the collation means 101, and the acoustic model managing server 20, The acoustic model adaptation function to perform the same processing as the acoustic model adaptation means 1004, The adaptation finishing acoustic model storing function to perform the same processing as the adaptation finishing acoustic model storing means 1005, It becomes a voice recognition program by the software which comprises an adaptation finishing acoustic model function preselection capability which performs the same processing as the adaptation finishing acoustic model selecting means 1006, and an acoustic model transmitting function which performs the same processing as the acoustic model transmitting means 110.

[0167]The voice recognition program recorded on a recording medium The software of the voice recognition equipment 10, The software of the acoustic model managing server 20 may be recorded on a separate recording medium, may be recorded on one recording medium, and it may transmit to the acoustic model managing server 20 or the voice recognition equipment 10, respectively from the voice recognition equipment 10 or the acoustic model managing server 20.

[0168]With as mentioned above, the sound language model managing server 20 which was connected to the network according to this Embodiment 5. By building an acoustic model [ finishing / adaptation ] with the voice data for adaptation based on a user's audio signal 100, and updating the acoustic model 102 of a user's voice recognition equipment 10 via a network. The effect that the recognition precision of speech recognition can be raised is acquired without applying a big burden to a user.

[0169]By what the acoustic model 102 is updated for via a network by the adaptation finishing acoustic model which carried out adaptation to the user according to this Embodiment 5. Even when a user uses several different collation means 101, the suitable acoustic model 102 for the utilization time of all the collation means 101 is used, and the effect that high recognition precision can be acquired is acquired.

[0170]

[Effect of the Invention]As mentioned above, the acoustic model in which the acoustic model managing server acquired and built the updated sound data according to this invention, By transmitting to voice recognition equipment via a network, and updating the acoustic model which voice recognition equipment refers to in the case of speech recognition by the acoustic model which the acoustic model managing server transmitted, It is effective in the ability to raise the recognition precision of speech recognition, without applying a big burden to a user.

[0171]According to this invention, it corresponds to the specific condition directed by ID from which the acoustic model managing server acquired and acquired ID which specifies the acoustic model which voice recognition equipment refers to in the case of speech recognition, Can raise the recognition precision of speech recognition, without applying a big burden to a user by reading the updated sound data, building the acoustic model depending on a specific condition, and transmitting to voice recognition equipment, and. By transmitting the acoustic model customized for specification via a network, even when a user uses several different voice recognition equipments, a suitable acoustic model is used, and it is effective in the ability to acquire high recognition precision.

[0172]According to this invention, a language model managing server the language model which acquired and built the updated language data, By transmitting to voice recognition equipment via a network, and updating the language model which voice recognition equipment refers to in the case of speech recognition by the language model which the language model managing server transmitted, It is effective in the ability to raise the recognition precision of speech recognition, without applying a big burden to a user.

[0173]According to this invention, it corresponds to the specific condition directed by ID from which the language model managing server acquired and acquired ID which specifies the language model which voice recognition equipment refers to in the case of speech recognition, Can raise the recognition precision of speech recognition, without applying a big burden to a user by reading the updated language data, building the language model depending on a specific condition, and transmitting to voice recognition equipment, and. By updating the language model customized for specification via a network, even when a user uses several different voice recognition equipments, a suitable language model is used, and it is effective in the ability to acquire high recognition precision.

[0174]The language data in which the language model managing server read the user dictionary, and was updated via the network with reference to the user dictionary with which voice recognition equipment registered the word on the occasion of speech recognition according to this invention, It is effective in the ability to raise the recognition precision of speech recognition, without applying a big burden to a user, even when a user dictionary becomes large by building the language model depending on a user dictionary with reference to the read user dictionary, and transmitting to voice recognition equipment.

[0175]The language data which according to this invention the language model managing server acquired the text which the user of voice recognition equipment used, and was updated, It is effective in the ability to raise the recognition precision of speech recognition, without applying a big burden to a user by building the language model depending on a text with reference to the acquired text, and transmitting to the above-mentioned voice recognition equipment.

[0176]ID as which voice recognition equipment specifies an acoustic model according to this invention, Acquire the voice data for adaptation from the

inputted audio signal, and the voice data acquired ID and for adaptation, Transmit to an acoustic model managing server via a network, and an acoustic model managing server, Using the transmitted voice data for adaptation, carry out adaptation of the initial acoustic model, and match and store an adaptation finishing acoustic model in ID to which it was transmitted, and. In response to the renewal instructions of an acoustic model from the outside, ID is received from voice recognition equipment via a network, The adaptation finishing acoustic model corresponding to ID which received is chosen and read out of the stored adaptation finishing acoustic model, By transmitting to voice recognition equipment via a network, and updating the acoustic model which voice recognition equipment refers to in the case of speech recognition by the adaptation finishing acoustic model which the acoustic model managing server transmitted, Can raise the recognition precision of speech recognition, without applying a big burden to a user, and. Even when a user uses several different voice recognition equipments by updating an acoustic model via a network by the adaptation finishing acoustic model which carried out adaptation to the user, a suitable acoustic model is used, It is effective in the ability to acquire high recognition precision.

[0177]According to this invention, from the acoustic model managing server to which voice recognition equipment was connected via the network. When the acoustic model built with the updated sound data was received and the collation means was provided with the acoustic model update means which updates the acoustic model referred to in the case of speech recognition by the acoustic model which received, It is effective in the ability to raise the recognition precision of speech recognition, without applying a big burden to a user.

[0178]According to this invention, an acoustic model update means from the acoustic model managing server connected via the network. The acoustic model depending on the specific condition of the acoustic model which a collation means refers to in the case of speech recognition built with the updated sound data is received, Can raise the recognition precision of speech recognition, without applying a big burden to a user, when a collation means updates the acoustic model referred to in the case of speech recognition by the acoustic model which received, and. By receiving the acoustic model customized for specification via a network, even when a user uses several different collation means, a suitable acoustic model is used, and it is effective in the ability to acquire high recognition precision.

[0179]According to this invention, from the language model managing server to which voice recognition equipment was connected via the network. When the language model built with the updated language data was received and the collation means was provided with the language model update means which updates the language model referred to in the case of speech recognition by the language model which received, It is effective in the ability to raise the recognition precision of speech recognition, without applying a big burden to a user.

[0180]According to this invention, a language model update means from the language model managing server connected via the network. The language model depending on the specific condition of the language model which a collation means refers to in the case of speech recognition built with the updated language data is received, Can raise the recognition precision of speech recognition, without applying a big burden to a user, when a collation means updates the language model referred to in the case of speech recognition by the language model which received, and. By receiving the language model customized for specification via a network, even when a user uses several different collation means, a suitable language model is used, and it is effective in the ability to acquire high recognition precision.

[0181]According to this invention, a collation means is provided with the user dictionary which registered the word referred to in the case of speech recognition, A language model update means from the language model managing server connected via the network. When the language model depending on a user dictionary built with the updated language data is received and a collation means updates the language model referred to in the case of speech recognition by the language model which received, It is effective in the ability to raise the recognition precision of speech recognition, without applying a big burden to a user, even when a user dictionary becomes large.

[0182]According to this invention, a language model update means from the language model managing server connected via the network. When the language model depending on the text which was built with the updated language data and which the user who performs speech recognition used is received and a collation means updates the language model referred to in the case of speech recognition by the language model which received, It is effective in the ability to raise the recognition precision of speech recognition, without applying a big burden to a user.

[0183]The acoustic model which asks for the probability of an audio acoustical observed value series according to this invention, The collation means which inputs an audio signal, performs speech recognition with reference to the above-mentioned acoustic model, and outputs a recognition result, The acoustic model ID acquiring means which acquires ID which specifies an acoustic model, and acquired ID are read, The voice acquisition means for adaptation which acquires the voice data for adaptation from the inputted audio signal, and transmits the voice data read ID and for the acquired adaptation to the acoustic model managing server to which it was connected via the network, From an acoustic model managing server, the adaptation finishing acoustic model by which adaptation was carried out with the voice data for adaptation corresponding to ID is received, When the collation means was provided with the acoustic model update means which updates the acoustic model referred to in the case of speech recognition by the adaptation finishing acoustic model which received, Can raise the recognition precision of speech recognition, without applying a big burden to a user, and. It is effective in the ability to use a suitable acoustic model and acquire high recognition precision by updating an acoustic model via a network, by the adaptation finishing acoustic model which carried out adaptation to the user, even when a user uses several different collation means.

[0184]The sound data acquisition means which acquires the updated sound data according to this invention, The acoustic model construction means which builds the acoustic model which reads the sound data updated in response to the renewal instructions of an acoustic model from the outside, and asks for the probability of an audio acoustical observed value series, It is effective in the ability to raise the recognition precision of speech recognition, without applying a big burden to a user by having had the acoustic model transmitting means which transmits the acoustic model built by the acoustic model construction means to the voice recognition equipment which performs speech recognition via a network.

[0185]The sound data acquisition means which acquires the updated sound data according to this invention, The updating acoustic model ID acquiring means which acquires ID which specifies the acoustic model referred to in the case of speech recognition in response to the renewal instructions of an acoustic model from the outside, The sound data reading means for the specification which reads the updated sound data corresponding to the specific condition directed by acquired ID, The acoustic model construction means for the specification which builds the acoustic model depending on a specific condition with reference to the read sound data which was updated, Can raise the recognition precision of speech recognition, without applying a big burden to a user by having had the acoustic model transmitting means which transmits the built acoustic model to voice recognition equipment via a network, and. By transmitting the acoustic model customized for specification via a network, even when a user uses several different voice recognition equipments, a suitable acoustic model is used, and it is effective in the ability to acquire high recognition precision.

[0186]The initial acoustic model before adaptation which asks for the probability of an audio acoustical observed value series according to this invention, The voice data for adaptation transmitted from the voice recognition equipment connected via the network, ID which specifies the acoustic model which voice recognition equipment refers to in the case of speech recognition is received, An acoustic model adaptation means to carry out adaptation of the initial acoustic model using the voice data for adaptation, to match an adaptation finishing acoustic model with ID which received, and to store in an adaptation finishing acoustic model storing means, In response to the renewal instructions of an acoustic model from the outside, ID is received from voice recognition equipment via a network, The adaptation finishing acoustic model selecting means which chooses and reads the adaptation finishing acoustic model corresponding to ID which received from an adaptation finishing acoustic model storing means, Can raise the recognition precision of speech recognition, without applying a big burden to a user by having had the acoustic model transmitting means which transmits the read adaptation finishing acoustic model to voice recognition equipment via a network, and. Even when a user uses several different voice recognition equipments by updating an acoustic model via a network by the adaptation finishing acoustic model which carried out adaptation to the user, a suitable acoustic model is used, It is effective in the ability to acquire high recognition precision.

[0187]The language data acquisition means which acquires the updated language data according to this invention, The language model construction means which builds the language model which reads the language data updated in response to the renewal instructions of a language model from the outside, and searches for the appearing probability of a word sequence, It is effective in the ability to raise the recognition precision of speech recognition, without applying a big burden to a user by having had the language model transmitting means which transmits the language model which the

language model construction means built to the voice recognition equipment which performs speech recognition via a network.

[0188]The language data acquisition means which acquires the updated language data according to this invention, The updating language model ID acquiring means which acquires ID which specifies the language model referred to in the case of speech recognition in response to the renewal instructions of a language model from the outside, The language data reading means for the specification which reads the updated language data corresponding to the specific condition directed by ID, The language model construction means for the specification which builds the language model depending on a specific condition with reference to the read language data which was updated, Can raise the recognition precision of speech recognition, without applying a big burden to a user by having had the language model transmitting means which transmits the built language model to voice recognition equipment via a network, and. By transmitting the language model customized for specification via a network, even when a user uses several different voice recognition equipments, a suitable language model is used, and it is effective in the ability to acquire high recognition precision.

[0189]The language data acquisition means which acquires the updated language data according to this invention, The user dictionary reading means which reads the user dictionary which the voice recognition equipment connected via the network refers to in response to the renewal instructions of a language model from the outside in the case of speech recognition, The user dictionary dependence language model construction means which builds the language model which read the updated language data and was dependent on the user dictionary, It is effective in the ability to raise the recognition precision of speech recognition, without applying a big burden to a user, even when a user dictionary becomes large by having had the language model transmitting means which transmits the built language model to voice recognition equipment via a network.

[0190]The language data acquisition means which acquires the updated language data according to this invention, The user use text acquisition means which acquires the text which the user used in response to the renewal instructions of a language model from the outside, The user use text dependence language model construction means which builds the language model which read the updated language data and depended for it on the acquired text, It is effective in the ability to raise the recognition precision of speech recognition, without applying a big burden to a user by having had the language model transmitting means which transmits the built language model to voice recognition equipment via a network.

[0191]The 1st step that acquires the updated sound data according to this invention, The 2nd step that reads the sound data updated in response to the renewal instructions of an acoustic model, and builds an acoustic model, By having had the 3rd step that transmits the built acoustic model via a network, and the 4th step updated by the acoustic model which received the acoustic model and received the acoustic model referred to in the case of speech recognition, It is effective in the ability to raise the recognition precision of speech recognition, without applying a big burden to a user.

[0192]The 1st step that acquires the updated language data according to this invention, The 2nd step that reads the language data updated in response to the renewal instructions of a language model, and builds a language model, By having had the 3rd step that transmits the built language model via a network, and the 4th step updated by the language model which received the language model and received the language model referred to in the case of speech recognition, It is effective in the ability to raise the recognition precision of speech recognition, without applying a big burden to a user.

[0193]The 1st step that acquires the updated sound data according to this invention, The 2nd step that acquires ID which specifies the acoustic model referred to in the case of speech recognition in response to the renewal instructions of an acoustic model, The 3rd step that reads the updated sound data corresponding to the specific condition directed by ID, The 4th step that builds the acoustic model depending on a specific condition with reference to the updated sound data, By having had the 5th step that transmits the built acoustic model via a network, and the 6th step updated by the acoustic model which received the acoustic model and received the acoustic model referred to in the case of speech recognition, Can raise the recognition precision of speech recognition, without applying a big burden to a user, and. By receiving the acoustic model customized for specification via a network, even when a user uses several different voice recognition methods, a suitable acoustic model is used, and it is effective in the ability to acquire high recognition precision.

[0194]The 1st step that acquires the updated language data according to this invention, The 2nd step that acquires ID which specifies the language model referred to in the case of speech recognition in response to the renewal instructions of a language model, The 3rd step that reads the updated language data corresponding to the specific condition directed by acquired ID, The 4th step that builds the language model depending on a specific condition with reference to the updated language data, By having had the 5th step that transmits the built language model via a network, and the 6th step updated by the language model which received the language model and received the language model referred to in the case of speech recognition, Can raise the recognition precision of speech recognition, without applying a big burden to a user, and. By receiving the language model customized for specification via a network, even when a user uses several different voice recognition methods, a suitable language model is used, and it is effective in the ability to acquire high recognition precision.

[0195]The 1st step that acquires the updated language data according to this invention, The 2nd step that reads the user dictionary referred to in the case of speech recognition in response to the renewal instructions of a language model, The 3rd step that builds the language model which read the updated language data and was dependent on the user dictionary, By having had the 4th step that transmits the built language model via a network, and the 5th step updated by the language model which received the language model and received the language model referred to in the case of speech recognition, It is effective in the ability to raise the recognition precision of speech recognition, without applying a big burden to a user, even when a user dictionary becomes large.

[0196]The 1st step that acquires the updated language data according to this invention, The 2nd step that acquires the text which the user used in response to the renewal instructions of a language model, The 3rd step that builds the language model which read the updated language data and was dependent on the text, By having had the 4th step that transmits the built language model via a network, and the 5th step updated by the language model which received the language model and received the language model referred to in the case of speech recognition, It is effective in the ability to raise the recognition precision of speech recognition, without applying a big burden to a user.

[0197]The 1st step that acquires ID which specifies an acoustic model according to this invention, Read acquired ID, acquire the voice data for adaptation from the inputted audio signal, and via a network, The 2nd step that transmits the voice data read ID and for the acquired adaptation, The 3rd step that carries out adaptation of the initial acoustic model before adaptation, and matches and stores an adaptation finishing acoustic model in ID which transmitted using the transmitted voice data for adaptation, In response to the renewal instructions of an acoustic model, ID acquired at the 1st step via the network is received, The 4th step that chooses and reads the adaptation finishing acoustic model corresponding to ID which received out of the adaptation finishing acoustic model stored at the 3rd step, The 5th step that transmits the adaptation finishing acoustic model read at the 4th step via a network, By having had the 6th step updated by the adaptation finishing acoustic model which received the adaptation finishing acoustic model which transmitted and received the acoustic model referred to in the case of speech recognition, Can raise the recognition precision of speech recognition, without applying a big burden to a user, and by the adaptation finishing acoustic model which carried out adaptation to the user. Even when a user uses several different voice recognition methods by updating an acoustic model via a network, a suitable acoustic model is used and it is effective in the ability to acquire high recognition precision.

[0198]The sound data acquisition function which acquires the updated sound data with the recording medium which recorded the voice recognition program according to this invention, The acoustic model construction function to read the sound data updated in response to the renewal instructions of an acoustic model, and to build an acoustic model, The acoustic model transmitting function which transmits the built acoustic model via a network, It is effective in the ability to raise the recognition precision of speech recognition, without applying a big burden to a user by receiving an acoustic model and realizing the acoustic model update function which a verification function updates by the acoustic model which received the acoustic model referred to in the case of speech recognition.

[0199]The language data acquisition function which acquires the updated language data with the recording medium which recorded the voice recognition program according to this invention, The language model construction function to read the updated language data which the above-mentioned language data acquisition function acquired in response to the renewal instructions of a language model, and to build a language model, The language model transmitting function which transmits the language model which the above-mentioned language model construction function built via a network, By



receiving the language model which the above-mentioned language model transmitting function transmitted, and realizing the language model update function which the above-mentioned verification function updates by the language model which received the language model referred to in the case of speech recognition, It is effective in the ability to raise the recognition precision of speech recognition, without applying a big burden to a user.

[0200]The sound data acquisition function which acquires the updated sound data with the recording medium which recorded the voice recognition program according to this invention, The updating acoustic model ID acquisition function which acquires ID which specifies an acoustic model in response to the renewal instructions of an acoustic model, The sound data read-out function for the specification which reads the updated sound data corresponding to the specific condition directed by acquired ID, The acoustic model construction function for the specification which builds the acoustic model depending on a specific condition with reference to the updated sound data, The acoustic model transmitting function which transmits an acoustic model via a network, By receiving an acoustic model and realizing the acoustic model update function which a verification function updates by the acoustic model which received the acoustic model referred to in the case of speech recognition, Can raise the recognition precision of speech recognition, without applying a big burden to a user, and. By receiving the acoustic model customized for specification via a network, even when a user uses two or more verification functions, a suitable acoustic model is used, and it is effective in the ability to acquire high recognition precision.

[0201]The language data acquisition function which acquires the updated language data with the recording medium which recorded the voice recognition program according to this invention, The updating language model ID acquisition function which acquires ID which specifies a language model in response to the renewal instructions of a language model, The language data read-out function for the specification which reads the updated language data corresponding to the specific condition directed by acquired ID, The language model construction function for the specification which builds the language model depending on a specific condition with reference to the updated language data, The language model transmitting function which transmits the built language model via a network, By receiving a language model and realizing the language model update function which a verification function updates by the language model which received the language model referred to in the case of speech recognition, Can raise the recognition precision of speech recognition, without applying a big burden to a user, and. By receiving the language model customized for specification via a network, even when a user uses two or more verification functions, a suitable language model is used, and it is effective in the ability to acquire high recognition precision.

[0202]The language data acquisition function which acquires the updated language data with the recording medium which recorded the voice recognition program according to this invention, The user dictionary read-out function which reads a user dictionary in response to the renewal instructions of a language model, The user dictionary dependence language model construction function to build the language model which read the updated language data and was dependent on the user dictionary, The language model transmitting function which transmits the built language model via a network, By receiving a language model and realizing the language model update function which a verification function updates by the language model which received the language model referred to in the case of speech recognition, It is effective in the ability to raise the recognition precision of speech recognition, without applying a big burden to a user, even when a user dictionary becomes large.

[0203]The language data acquisition function which acquires the updated language data with the recording medium which recorded the voice recognition program according to this invention, The user use text acquisition function which acquires the text which the user used in response to the renewal instructions of a language model, The user use text dependence language model construction function to build the language model which read the updated language data and was dependent on the text, The language model transmitting function which transmits the built language model via a network, It is effective in the ability to raise the recognition precision of speech recognition, without applying a big burden to a user by receiving a language model and realizing the language model update function which a verification function updates by the language model which received the language model referred to in the case of speech recognition.

[0204]The acoustic model ID acquisition function which acquires ID which specifies an acoustic model with the recording medium which recorded the voice recognition program according to this invention, Read acquired ID, acquire the voice data for adaptation from the inputted audio signal, and via a network, The voice acquisition function for adaptation which transmits the voice data read ID and for the acquired adaptation, The acoustic model adaptation function to carry out adaptation of the initial acoustic model before adaptation, and to match and store an adaptation finishing acoustic model in ID which transmitted using the transmitted voice data for adaptation, In response to the renewal instructions of an acoustic model, ID which the acoustic model ID acquisition function acquired via the network is received, The adaptation finishing acoustic model function preselection capability which chooses and reads the adaptation finishing acoustic model corresponding to ID which received out of the adaptation finishing acoustic model which the acoustic model adaptation function stored, The acoustic model transmitting function which transmits the read adaptation finishing acoustic model via a network, Without applying a big burden to a user by receiving an adaptation finishing acoustic model and realizing the acoustic model update function which a verification function updates by the adaptation finishing acoustic model which received the acoustic model referred to in the case of speech recognition, Can raise the recognition precision of speech recognition and by updating an acoustic model via a network by the adaptation finishing acoustic model which carried out adaptation to the user. Even when a user uses several different verification functions, a suitable acoustic model is used and it is effective in the ability to acquire high recognition precision.

---

[Translation done.]

## \* NOTICES \*

JPO and INPIT are not responsible for any damages caused by the use of this translation.

- 1.This document has been translated by computer. So the translation may not reflect the original precisely.
- 2.\*\*\* shows the word which can not be translated.
- 3.In the drawings, any words are not translated.

---

DESCRIPTION OF DRAWINGS

---

## [Brief Description of the Drawings]

- [Drawing 1] It is a block diagram showing the composition of the voice recognition equipment by this embodiment of the invention 1.
- [Drawing 2] It is a flow chart which shows the procedure of the update process of the acoustic model by this embodiment of the invention 1.
- [Drawing 3] It is a flow chart which shows the procedure of the update process of the language model by this embodiment of the invention 1.
- [Drawing 4] It is a block diagram showing the composition of the voice recognition equipment by this embodiment of the invention 2.
- [Drawing 5] It is a flow chart which shows the procedure of the update process of the language model by this embodiment of the invention 2.
- [Drawing 6] It is a block diagram showing the composition of the voice recognition equipment by this embodiment of the invention 3.
- [Drawing 7] It is a flow chart which shows the procedure of the update process of the language model by this embodiment of the invention 3.
- [Drawing 8] It is a block diagram showing the composition of the voice recognition equipment by this embodiment of the invention 4.
- [Drawing 9] It is a flow chart which shows the procedure of the update process of the language model by this embodiment of the invention 4.
- [Drawing 10] It is a block diagram showing the composition of the voice recognition equipment by this embodiment of the invention 5.
- [Drawing 11] It is a flow chart which shows the procedure of the update process of the acoustic model by this embodiment of the invention 5.
- [Drawing 12] It is a block diagram showing the composition of conventional voice recognition equipment.
- [Drawing 13] It is a flow chart which shows the procedure of the conventional speech recognition processing.
- [Drawing 14] It is a block diagram showing the composition of conventional voice recognition equipment.
- [Drawing 15] It is a block diagram showing the composition of conventional voice recognition equipment.
- [Drawing 16] It is a figure showing the example of composition of the user dictionary 601.
- [Drawing 17] It is a flow chart which shows the procedure of the conventional speech recognition processing.

## [Description of Notations]

10 Voice recognition equipment, 20 acoustic-model managing server, 30 language-model managing server, 100 audio signals and 101 A collation means, 102 acoustic models, 103 language models, 104 A recognition result, the renewal instructions of 105 acoustic models, a 106 sound-data acquisition means, 107 The updated sound data, a 108 acoustic-model construction means, a 109 acoustic-model storing means, 110 An acoustic model transmitting means, a 111 acoustic-model update means, the renewal instructions of 112 language models, 113 A language data acquisition means, the language data updated 114 times, a 115 language-model construction means, 116 A language model storing means, a 117 language-model transmitting means, a 118 language-model update means, The renewal language model ID acquiring means of 401, and 402 The language data reading means for specification, 403 The language model construction means for specification, 601 user dictionaries, a 602 user-dictionary reading means, 603 A user dictionary dependence language model construction means and 801 User use text acquisition means, 802 A user use text storing means and 803 A user use text dependence language model construction means, a 1001 acoustic-model ID acquiring means, and 1002 The voice acquisition means for adaptation, 1003 initial acoustic models, 1004 An acoustic model adaptation means and 1005 An adaptation finishing acoustic model storing means and 1006 Adaptation finishing acoustic model selecting means.

---

[Translation done.]



(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開2002-91477

(P2002-91477A)

(43) 公開日 平成14年3月27日 (2002.3.27)

(51) Int.Cl.<sup>7</sup>

識別記号

F I

テーマコード\* (参考)

G 1 0 L 15/06  
15/14  
15/18  
15/28

G 1 0 L 3/00

5 2 1 A 5 D 0 1 5

5 2 1 R

5 3 5 Z

5 3 7 D

5 7 1 A

審査請求 未請求 請求項の数35 O L (全 32 頁)

(21) 出願番号 特願2000-280674(P2000-280674)

(22) 出願日 平成12年9月14日 (2000.9.14)

(71) 出願人 000006013

三菱電機株式会社

東京都千代田区丸の内二丁目2番3号

(72) 発明者 岡登 洋平

東京都千代田区丸の内二丁目2番3号 三  
菱電機株式会社内

(72) 発明者 石井 純

東京都千代田区丸の内二丁目2番3号 三  
菱電機株式会社内

(74) 代理人 100066474

弁理士 田澤 博昭 (外1名)

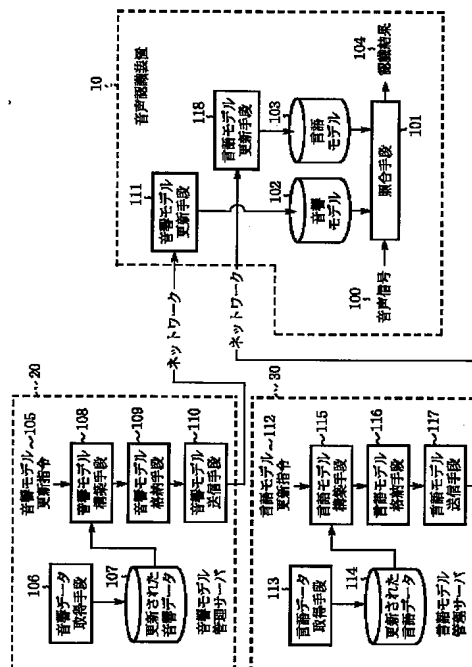
Fターム(参考) 5D015 GG05 HH05 HH23

(54) 【発明の名称】 音声認識システム、音声認識装置、音響モデル管理サーバ、言語モデル管理サーバ、音声認識方法及び音声認識プログラムを記録したコンピュータ読み取り可能な記録媒体

(57) 【要約】

【課題】 ユーザへ大きな負担をかけることなく音響モデルと言語モデルを更新し、認識精度を向上させる。

【解決手段】 音響モデル管理サーバ20は、更新された音響データ107を取得し音響モデルを構築し、言語モデル管理サーバ30は、更新された言語データ114を取得し言語モデルを構築して、それぞれネットワークを介して音声認識装置10に送信する。音声認識装置10において、音響モデル更新手段111は、送信された音響モデルにより音響モデル102を更新し、言語モデル更新手段118は、送信された言語モデルにより言語モデル103を更新する。



## 【特許請求の範囲】

【請求項1】 音声信号を入力し、音声の音響的な観測値系列の確率を求める音響モデルを参照して音声認識を行い、認識結果を出力する音声認識装置と、

上記音声認識装置とネットワークを介して接続され、更新された音響データを取得して上記音響モデルを構築する音響モデル管理サーバとを備えた音声認識システムにおいて、

上記音響モデル管理サーバが構築した上記音響モデルを上記音声認識装置に送信し、上記音声認識装置が、音声認識の際に参照する音響モデルを、上記音響モデル管理サーバが送信した音響モデルにより更新することを特徴とする音声認識システム。

【請求項2】 音響モデル管理サーバが、音声認識装置が音声認識の際に参照する音響モデルを特定するIDを取得し、取得したIDで指示される特定条件に対応して、更新された音響データを読み出し、上記特定条件に依存した音響モデルを構築して上記音声認識装置に送信することを特徴とする請求項1記載の音声認識システム。

【請求項3】 音声信号を入力し、単語列の出現確率を求める言語モデルを参照して音声認識を行い、認識結果を出力する音声認識装置と、

上記音声認識装置とネットワークを介して接続され、更新された言語データを取得して上記言語モデルを構築する言語モデル管理サーバとを備えた音声認識システムにおいて、

上記言語モデル管理サーバが構築した上記言語モデルを上記音声認識装置に送信し、上記音声認識装置が、音声認識の際に参照する言語モデルを、上記言語モデル管理サーバが送信した言語モデルにより更新することを特徴とする音声認識システム。

【請求項4】 言語モデル管理サーバが、音声認識装置が音声認識の際に参照する言語モデルを特定するIDを取得し、取得したIDで指示される特定条件に対応して、更新された言語データを読み出し、上記特定条件に依存した言語モデルを構築して上記音声認識装置に送信することを特徴とする請求項3記載の音声認識システム。

【請求項5】 音声認識装置が音声認識の際に単語を登録したユーザ辞書を参照し、

言語モデル管理サーバが、ネットワークを介して上記ユーザ辞書を読み出し、更新された言語データと、読み出した上記ユーザ辞書とを参照し、上記ユーザ辞書に依存した言語モデルを構築して上記音声認識装置に送信することを特徴とする請求項3記載の音声認識システム。

【請求項6】 言語モデル管理サーバが、音声認識装置のユーザが利用したテキストを取得し、更新された言語データと、取得した上記テキストとを参照し、上記テキストに依存した言語モデルを構築して上記音声認識装置

に送信することを特徴とする請求項3記載の音声認識システム。

【請求項7】 音声信号を入力し、音声の音響的な観測値系列の確率を求める音響モデルを参照して音声認識を行い、認識結果を出力する音声認識装置と、

上記音声認識装置とネットワークを介して接続され、適応化前の初期音響モデルを有する音響モデル管理サーバとを備えた音声認識システムにおいて、

上記音声認識装置が、上記音響モデルを特定するIDと、入力された音声信号から適応化用の音声データとを取得し、取得したID及び適応化用の音声データを、ネットワークを介して上記音響モデル管理サーバに送信し、

上記音響モデル管理サーバが、送信された適応化用の音声データを用いて、上記初期音響モデルを適応化し、適応化済み音響モデルを、送信された上記IDに対応付けて格納すると共に、外部からの音響モデル更新指令を受けて、ネットワークを介して上記音声認識装置から上記音響モデルを特定するIDを受信し、受信したIDに対応する適応化済み音響モデルを、格納している適応化済み音響モデルの中から選択して読み出し、ネットワークを介して上記音声認識装置に送信し、

上記音声認識装置が、音声認識の際に参照する音響モデルを、上記音響モデル管理サーバが送信した適応化済み音響モデルにより更新することを特徴とする音声認識システム。

【請求項8】 音声の音響的な観測値系列の確率を求める音響モデルと、

音声信号を入力し上記音響モデルを参照して音声認識を行い、認識結果を出力する照合手段とを備えた音声認識装置において、

ネットワークを介して接続された音響モデル管理サーバから、更新された音響データにより構築された音響モデルを受信し、上記照合手段が音声認識の際に参照する音響モデルを、受信した音響モデルにより更新する音響モデル更新手段とを備えたことを特徴とする音声認識装置。

【請求項9】 音響モデル更新手段が、ネットワークを介して接続された音響モデル管理サーバから、更新された音響データにより構築された、照合手段が音声認識の際に参照する音響モデルの特定条件に依存した音響モデルを受信し、上記照合手段が音声認識の際に参照する音響モデルを、受信した音響モデルにより更新することを特徴とする請求項8記載の音声認識装置。

【請求項10】 単語列の出現確率を求める言語モデルと、音声信号を入力し上記言語モデルを参照して音声認識を行い、認識結果を出力する照合手段とを備えた音声認識装置において、

ネットワークを介して接続された言語モデル管理サーバ

から、更新された言語データにより構築された言語モデルを受信し、上記照合手段が音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する言語モデル更新手段とを備えたことを特徴とする音声認識装置。

【請求項11】 言語モデル更新手段が、ネットワークを介して接続された言語モデル管理サーバから、更新された言語データにより構築された、照合手段が音声認識の際に参照する言語モデルの特定条件に依存した言語モデルを受信し、上記照合手段が音声認識の際に参照する言語モデルを、受信した言語モデルにより更新することを特徴とする請求項10記載の音声認識装置。

【請求項12】 照合手段が音声認識の際に参照する単語を登録したユーザ辞書を備え、言語モデル更新手段が、ネットワークを介して接続された言語モデル管理サーバから、更新された言語データにより構築された、上記照合手段が音声認識の際に参照するユーザ辞書に依存した言語モデルを受信し、上記照合手段が音声認識の際に参照する言語モデルを、受信した言語モデルにより更新することを特徴とする請求項10記載の音声認識装置。

【請求項13】 言語モデル更新手段が、ネットワークを介して接続された言語モデル管理サーバから、更新された言語データにより構築された、音声認識を行うユーザが利用したテキストに依存した言語モデルを受信し、上記照合手段が音声認識の際に参照する言語モデルを、受信した言語モデルにより更新することを特徴とする請求項10記載の音声認識装置。

【請求項14】 音声の音響的な観測値系列の確率を求める音響モデルと、音声信号を入力し上記音響モデルを参照して音声認識を行い、認識結果を出力する照合手段と、上記音響モデルを特定するIDを取得する音響モデルID取得手段と、上記音響モデルID取得手段が取得したIDを読み出し、入力された音声信号から適応化用の音声データを取得し、読み出したID及び取得した適応化用の音声データを、ネットワークを介して接続された音響モデル管理サーバに送信する適応化用音声取得手段と、上記音響モデル管理サーバから、上記IDに対応する上記適応化用の音声データにより適応化された適応化済み音響モデルを受信し、上記照合手段が音声認識の際に参照する音響モデルを、受信した適応化済み音響モデルにより更新する音響モデル更新手段とを備えたことを特徴とする音声認識装置。

【請求項15】 更新された音響データを取得する音響データ取得手段と、外部からの音響モデル更新指令を受けて、上記音響データ取得手段が取得した更新された音響データを読み出し、音声の音響的な観測値系列の確率を求める音響モ

デルを構築する音響モデル構築手段と、上記音響モデル構築手段により構築された音響モデルを、ネットワークを介して音声認識を行う音声認識装置に送信する音響モデル送信手段とを備えたことを特徴とする音響モデル管理サーバ。

【請求項16】 更新された音響データを取得する音響データ取得手段と、

外部からの音響モデル更新指令を受けて、ネットワークを介して接続された音声認識装置が音声認識の際に参照する音響モデルを特定するIDを取得する更新音響モデルID取得手段と、

上記更新音響モデルID取得手段が取得したIDで指示される特定条件に対応して、上記音響データ取得手段が取得した更新された音響データを読み出す特定向け音響データ読み出し手段と、

上記特定向け音響データ読み出し手段が読み出した更新された音響データを参照し、上記特定条件に依存した音響モデルを構築する特定向け音響モデル構築手段と、

上記特定向け音響モデル構築手段が構築した音響モデルを、ネットワークを介して上記音声認識装置に送信する音響モデル送信手段とを備えたことを特徴とする音響モデル管理サーバ。

【請求項17】 音声の音響的な観測値系列の確率を求める、適応化前の初期音響モデルと、

ネットワークを介して接続された音声認識装置から送信された、適応化用の音声データと、上記音声認識装置が音声認識の際に参照する音響モデルを特定するIDを受信し、上記適応化用の音声データを用いて上記初期音響モデルを適応化し、適応化済み音響モデルを、受信した

上記IDに対応付けて適応化済み音響モデル格納手段に格納する音響モデル適応化手段と、

外部からの音響モデル更新指令を受けて、ネットワークを介して上記音声認識装置から上記IDを受信し、受信したIDに対応する適応化済み音響モデルを、上記適応化済み音響モデル格納手段から選択して読み出す適応化済み音響モデル選択手段と、

上記適応化済み音響モデル選択手段が読み出した適応化済み音響モデルを、ネットワークを介して上記音声認識装置に送信する音響モデル送信手段とを備えたことを特徴とする音響モデル管理サーバ。

【請求項18】 更新された言語データを取得する言語データ取得手段と、

外部からの言語モデル更新指令を受けて、上記言語データ取得手段が取得した更新された言語データを読み出し、単語列の出現確率を求める言語モデルを構築する言語モデル構築手段と、

上記言語モデル構築手段が構築した言語モデルを、ネットワークを介して音声認識を行う音声認識装置に送信する言語モデル送信手段とを備えたことを特徴とする言語モデル管理サーバ。

【請求項19】 更新された言語データを取得する言語データ取得手段と、

外部からの言語モデル更新指令を受けて、ネットワークを介して接続された音声認識装置が音声認識の際に参照する言語モデルを特定するIDを取得する更新言語モデルID取得手段と、

上記更新言語モデルID取得手段が取得したIDで指示される特定条件に対応して、上記言語データ取得手段が取得した更新された言語データを読み出す特定向け言語データ読み出し手段と、

上記特定向け言語データ読み出し手段が読み出した更新された言語データを参照し、上記特定条件に依存した言語モデルを構築する特定向け言語モデル構築手段と、

上記特定向け言語モデル構築手段が構築した言語モデルを、ネットワークを介して上記音声認識装置に送信する言語モデル送信手段とを備えたことを特徴とする言語モデル管理サーバ。

【請求項20】 更新された言語データを取得する言語データ取得手段と、

外部からの言語モデル更新指令を受けて、ネットワークを介して接続された音声認識装置が音声認識の際に参照するユーザ辞書を読み出すユーザ辞書読み出し手段と、上記言語データ取得手段が取得した更新された言語データを読み出し、上記ユーザ辞書読み出し手段が読み出したユーザ辞書に依存した言語モデルを構築するユーザ辞書依存言語モデル構築手段と、

上記ユーザ辞書依存言語モデル構築手段が構築した言語モデルを、ネットワークを介して上記音声認識装置に送信する言語モデル送信手段とを備えたことを特徴とする言語モデル管理サーバ。

【請求項21】 更新された言語データを取得する言語データ取得手段と、

外部からの言語モデル更新指令を受けて、ネットワークを介して接続された音声認識装置のユーザが利用したテキストを取得するユーザ利用テキスト取得手段と、

上記言語データ取得手段が取得した更新された言語データを読み出し、上記ユーザ利用テキスト取得手段が取得したテキストに依存した言語モデルを構築するユーザ利用テキスト依存言語モデル構築手段と、

上記ユーザ利用テキスト依存言語モデル構築手段が構築した言語モデルを、ネットワークを介して上記音声認識装置に送信する言語モデル送信手段とを備えたことを特徴とする言語モデル管理サーバ。

【請求項22】 音声信号を入力し、音声の音響的な観測値系列の確率を求める音響モデルを参照して音声認識を行い、認識結果を出力する音声認識方法において、更新された音響データを取得する第1のステップと、音響モデル更新指令を受けて、上記第1のステップで取得した更新された音響データを読み出し、音響モデルを構築する第2のステップと、

上記第2のステップで構築した音響モデルを、ネットワークを介して送信する第3のステップと、

上記第3のステップで送信した音響モデルを受信し、上記音声認識の際に参照する音響モデルを、受信した音響モデルにより更新する第4のステップとを備えたことを特徴とする音声認識方法。

【請求項23】 音声信号を入力し、単語列の出現確率を求める言語モデルを参照して音声認識を行い、認識結果を出力する音声認識方法において、

10 更新された言語データを取得する第1のステップと、

言語モデル更新指令を受けて、上記第1のステップで取得した更新された言語データを読み出し、言語モデルを構築する第2のステップと、

上記第2のステップで構築した言語モデルを、ネットワークを介して送信する第3のステップと、

上記第3のステップで送信した言語モデルを受信し、上記音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する第4のステップとを備えたことを特徴とする音声認識方法。

20 【請求項24】 音声信号を入力し、音声の音響的な観測値系列の確率を求める音響モデルを参照して音声認識を行い、認識結果を出力する音声認識方法において、更新された音響データを取得する第1のステップと、音響モデル更新指令を受けて、音声認識の際に参照する音響モデルを特定するIDを取得する第2のステップと、

上記第2のステップで取得したIDで指示される特定条件に対応して、上記第1のステップで取得した更新された音響データを読み出す第3のステップと、

30 上記第3のステップで読み出した更新された音響データを参照し、上記特定条件に依存した音響モデルを構築する第4のステップと、

上記第4のステップで構築した音響モデルを、ネットワークを介して送信する第5のステップと、

上記第5のステップで送信した音響モデルを受信し、音声認識の際に参照する音響モデルを、受信した音響モデルにより更新する第6のステップとを備えたことを特徴とする音声認識方法。

40 【請求項25】 音声信号を入力し、単語列の出現確率を求める言語モデルを参照して音声認識を行い、認識結果を出力する音声認識方法において、

更新された言語データを取得する第1のステップと、

言語モデル更新指令を受けて、音声認識の際に参照する言語モデルを特定するIDを取得する第2のステップと、

上記第2のステップで取得したIDで指示される特定条件に対応して、上記第1のステップで取得した更新された言語データを読み出す第3のステップと、

50 上記第3のステップで読み出した更新された言語データを参照し、上記特定条件に依存した言語モデルを構築す

る第4のステップと、  
上記第4のステップで構築した言語モデルを、ネットワークを介して送信する第5のステップと、  
上記第5のステップで送信した言語モデルを受信し、音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する第6のステップとを備えたことを特徴とする音声認識方法。

【請求項26】 音声信号を入力し、単語列の出現確率を求める言語モデルと、単語を登録したユーザ辞書を参照して音声認識を行い、認識結果を出力する音声認識方法において、

更新された言語データを取得する第1のステップと、  
言語モデル更新指令を受けて、音声認識の際に参照するユーザ辞書を読み出す第2のステップと、  
上記第1のステップで取得した更新された言語データを読み出し、上記第2のステップで読み出したユーザ辞書に依存した言語モデルを構築する第3のステップと、  
上記第3のステップで構築した言語モデルを、ネットワークを介して送信する第4のステップと、  
上記第4のステップで送信した言語モデルを受信し、音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する第5のステップとを備えたことを特徴とする音声認識方法。

【請求項27】 音声信号を入力し、単語列の出現確率を求める言語モデルを参照して音声認識を行い、認識結果を出力する音声認識方法において、  
更新された言語データを取得する第1のステップと、  
言語モデル更新指令を受けて、音声認識を行うユーザが利用したテキストを取得する第2のステップと、  
上記第1のステップで取得した更新された言語データを読み出し、上記第2のステップで取得したテキストに依存した言語モデルを構築する第3のステップと、  
上記第3のステップで構築した言語モデルを、ネットワークを介して送信する第4のステップと、  
上記第4のステップで送信した言語モデルを受信し、音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する第5のステップとを備えたことを特徴とする音声認識方法。

【請求項28】 音声信号を入力し、音声の音響的な観測値系列の確率を求める音響モデルを参照して音声認識を行い、認識結果を出力する音声認識方法において、  
上記音響モデルを特定するIDを取得する第1のステップと、  
上記第1のステップで取得したIDを読み出し、入力された音声信号から適応化用の音声データを取得し、ネットワークを介して、読み出したID及び取得した適応化用の音声データを送信する第2のステップと、  
上記第2のステップで送信した適応化用の音声データを用いて、適応化前の初期音響モデルを適応化し、適応化済み音響モデルを、上記第2のステップで送信したID

に対応付けて格納する第3のステップと、  
音響モデル更新指令を受けて、ネットワークを介して上記第1のステップで取得したIDを受信し、受信したIDに対応する適応化済み音響モデルを、上記第3のステップで格納している適応化済み音響モデルの中から選択して読み出す第4のステップと、  
上記第4のステップで読み出した適応化済み音響モデルを、ネットワークを介して送信する第5のステップと、  
上記第5のステップで送信した適応化済み音響モデルを受信し、音声認識の際に参照する音響モデルを、受信した適応化済み音響モデルにより更新する第6のステップとを備えたことを特徴とする音声認識方法。

【請求項29】 音声信号を入力し、音声の音響的な観測値系列の確率を求める音響モデルを参照して音声認識を行い、認識結果を出力する照合機能を実現させる音声認識プログラムを記録した記録媒体であって、  
更新された音響データを取得する音響データ取得機能と、  
音響モデル更新指令を受けて、上記音響データ取得機能が取得した更新された音響データを読み出し、音響モデルを構築する音響モデル構築機能と、  
上記音響モデル構築機能が構築した音響モデルを、ネットワークを介して送信する音響モデル送信機能と、  
上記音響モデル送信機能が送信した音響モデルを受信し、上記照合機能が音声認識の際に参照する音響モデルを、受信した音響モデルにより更新する音響モデル更新機能とを実現させる音声認識プログラムを記録したコンピュータ読み取り可能な記録媒体。

【請求項30】 音声信号を入力し、単語列の出現確率を求める言語モデルを参照して音声認識を行い、認識結果を出力する照合機能を実現させる音声認識プログラムを記録した記録媒体であって、  
更新された言語データを取得する言語データ取得機能と、  
言語モデル更新指令を受けて、上記言語データ取得機能が取得した更新された言語データを読み出し、言語モデルを構築する言語モデル構築機能と、  
上記言語モデル構築機能が構築した言語モデルを、ネットワークを介して送信する言語モデル送信機能と、  
上記言語モデル送信機能が送信した言語モデルを受信し、上記照合機能が音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する言語モデル更新機能とを実現させる音声認識プログラムを記録したコンピュータ読み取り可能な記録媒体。

【請求項31】 音声信号を入力し、音声の音響的な観測値系列の確率を求める音響モデルを参照して音声認識を行い、認識結果を出力する照合機能を実現させる音声認識プログラムを記録した記録媒体であって、  
更新された音響データを取得する音響データ取得機能と、音響モデル更新指令を受けて、上記音響モデルを特

定するIDを取得する更新音響モデルID取得機能と、  
上記更新音響モデルID取得機能が取得したIDで指示  
される特定条件に対応して、上記音響データ取得機能が  
取得した更新された音響データを読み出す特定向け音響  
データ読み出し機能と、

上記特定向け音響データ読み出し機能が読み出した更新  
された音響データを参照し、上記特定条件に依存した音  
響モデルを構築する特定向け音響モデル構築機能と、  
上記特定向け音響モデル構築機能が構築した音響モデル  
を、ネットワークを介して送信する音響モデル送信機能  
と、

上記音響モデル送信機能が送信した音響モデルを受信  
し、上記照合機能が音声認識の際に参照する音響モデル  
を、受信した音響モデルにより更新する音響モデル更新  
機能とを実現させる音声認識プログラムを記録したコン  
ピュータ読み取り可能な記録媒体。

【請求項32】 音声信号を入力し、単語列の出現確率  
を求める言語モデルを参照して音声認識を行い、認識結  
果を出力する照合機能を実現させる音声認識プログラム  
を記録した記録媒体であって、  
更新された言語データを取得する言語データ取得機能  
と、

言語モデル更新指令を受けて、上記言語モデルを特定す  
るIDを取得する更新言語モデルID取得機能と、  
上記更新言語モデルID取得機能が取得したIDで指示  
される特定条件に対応して、上記言語データ取得機能が  
取得した更新された言語データを読み出す特定向け言語  
データ読み出し機能と、

上記特定向け言語データ読み出し機能が読み出した更新  
された言語データを参照し、上記特定条件に依存した言  
語モデルを構築する特定向け言語モデル構築機能と、  
上記特定向け言語モデル構築機能が構築した言語モデル  
を、ネットワークを介して送信する言語モデル送信機能  
と、

上記言語モデル送信機能が送信した言語モデルを受信  
し、上記照合機能が音声認識の際に参照する言語モデル  
を、受信した言語モデルにより更新する言語モデル更新  
機能とを実現させる音声認識プログラムを記録したコン  
ピュータ読み取り可能な記録媒体。

【請求項33】 音声信号を入力し、単語列の出現確率  
を求める言語モデルと、単語を登録したユーザ辞書を参  
照して音声認識を行い、認識結果を出力する照合機能  
を実現させる音声認識プログラムを記録した記録媒体  
であって、

更新された言語データを取得する言語データ取得機能  
と、

言語モデル更新指令を受けて、上記ユーザ辞書を読み出  
すユーザ辞書読み出し機能と、

上記言語データ取得機能が取得した更新された言語デ  
ータを読み出し、上記ユーザ辞書読み出し機能が読み出し

たユーザ辞書に依存した言語モデルを構築するユーザ辞  
書依存言語モデル構築機能と、

上記ユーザ辞書依存言語モデル構築機能が構築した言語  
モデルを、ネットワークを介して送信する言語モデル送  
信機能と、

上記言語モデル送信機能が送信した言語モデルを受信  
し、上記照合機能が音声認識の際に参照する言語モデル  
を、受信した言語モデルにより更新する言語モデル更新  
機能とを実現させる音声認識プログラムを記録したコン  
ピュータ読み取り可能な記録媒体。

【請求項34】 音声信号を入力し、単語列の出現確率  
を求める言語モデルを参照して音声認識を行い、認識結  
果を出力する照合機能を実現させる音声認識プログラム  
を記録した記録媒体であって、

更新された言語データを取得する言語データ取得機能  
と、

言語モデル更新指令を受けて、音声認識を行うユーザが  
利用したテキストを取得するユーザ利用テキスト取得機  
能と、

上記言語データ取得機能が取得した更新された言語デ  
ータを読み出し、上記ユーザ利用テキスト取得機能が取得  
したテキストに依存した言語モデルを構築するユーザ利  
用テキスト依存言語モデル構築機能と、

上記ユーザ利用テキスト依存言語モデル構築機能が構築  
した言語モデルを、ネットワークを介して送信する言語  
モデル送信機能と、

上記言語モデル送信機能が送信した言語モデルを受信  
し、上記照合機能が音声認識の際に参照する言語モデル  
を、受信した言語モデルにより更新する言語モデル更新  
機能とを実現させる音声認識プログラムを記録したコン  
ピュータ読み取り可能な記録媒体。

【請求項35】 音声信号を入力し、音声の音響的な観  
測値系列の確率を求める音響モデルを参照して音声認識  
を行い、認識結果を出力する照合機能を実現させる音声  
認識プログラムを記録した記録媒体であって、

上記音響モデルを特定するIDを取得する音響モデルID  
取得機能と、

上記音響モデルID取得機能が取得したIDを読み出  
し、入力された音声信号から適応化用の音声データを取  
得し、ネットワークを介して、読み出したID及び取得  
した適応化用の音声データを送信する適応化用音声取得  
機能と、

上記適応化用音声取得機能が送信した適応化用の音声デ  
ータを用いて、適応化前の初期音響モデルを適応化し、  
適応化済み音響モデルを、上記適応化用音声取得機能が  
送信したIDに対応付けて格納する音響モデル適応化機  
能と、

音響モデル更新指令を受けて、ネットワークを介して上  
記音響モデルID取得機能が取得したIDを受信し、受  
信したIDに対応する適応化済み音響モデルを、上記音

響モデル適応化機能が格納した適応化済み音響モデルの中から選択して読み出す適応化済み音響モデル選択機能と、

上記適応化済み音響モデル選択機能が読み出した適応化済み音響モデルを、ネットワークを介して送信する音響モデル送信機能と、

上記音響モデル送信機能が送信した適応化済み音響モデルを受信し、上記照合機能が音声認識の際に参照する音響モデルを、受信した適応化済み音響モデルにより更新する音響モデル更新機能とを実現させる音声認識プログラムを記録したコンピュータ読み取り可能な記録媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】この発明は、音声認識の際に参照する音響モデルと言語モデルを、高い認識率が得られるように、ネットワークを介して最新の状態に更新す\*

$$W^* = \arg \max W P(W|A)$$

ただし、 $P(W|A)$ を直接求めることは通常困難である。そこでベイズの定理を用いて、 $P(W|A)$ は ※

$$P(W|A) = P(W) P(A|W) / P(A) \quad (2)$$

【0004】ここで、左辺を最大化する $W$ を求める際に、右辺分母である $P(A)$ は認識候補となる $W$ に影響★

$$W^* = \arg \max W P(W) P(A|W) \quad (3)$$

ここで、 $P(W)$ を与える確率モデルを言語モデルと呼び、 $P(A|W)$ を与える確率モデルを音響モデルと呼ぶ。

【0005】音声認識におけるこれらの代表的なモデル化方法は、音響モデルを隠れマルコフモデルで表現し、言語モデルを $n$ グラムと呼ばれる単語の $n-1$ 重マルコフ過程で表現する方法である。

【0006】これらの方法の詳細は、例えば「音声認識の基礎(上、下)」L. RABINER, B. H. JUANG, 古井監訳、1995年、11月、NTTアドバンステクノロジー(以下、文献1とする)、「確率的言語モデル」北研二、東京大学出版会(以下、文献2とする)、「音声・音情報のデジタル信号処理」鹿野清宏、中村哲、伊勢史郎共著、1997年、11月、昭晃堂(以下、文献3とする)に記されている。

【0007】これらの方法において、モデルを構成するパラメータは大量のデータから統計的に推定する。音響☆40

$$\begin{aligned} P(w_1 \dots w_k) \\ &= P(w_1) P(w_2|w_1) \dots P(w_k|w_1 \dots w_{k-1}) \\ &\approx P(w_1) P(w_2|w_1) \dots P(w_k|w_{k-1}) \quad (4) \end{aligned}$$

【0009】このように大量テキストを用いて、それぞれの単語の出現確率を統計的に推定することによって、統計量を用いない方法に比べて高い認識精度を得る言語モデルを構築できる。なお、日本語ではテキストが分かち書きされないため、単語の定義はあいまいであるが、本文では整合性のある何らかの手段でテキストを分割したそれぞれの単位を単語と定義する。この単語は、例え

＊る音声認識システム、音声認識装置、音響モデル管理サーバ、言語モデル管理サーバ、音声認識方法及び音声認識プログラムを記録したコンピュータ読み取り可能な記録媒体に関するものである。

【0002】

【従来の技術】音声認識においては、通常、デジタル化して入力された音声を、信号処理手法を用いて音声の音響的特徴を良く表すベクトルの時系列へ変換した後、音声のモデル(音響モデル、言語モデル)との照合処理を行う。

【0003】照合処理は $n$ 個の時刻フレームからなる音響特徴ベクトル時系列 $A = [a_1 a_2 : : a_n]$ から発生された単語列 $W = [w_1 w_2 : : w_k]$ ( $k$ は単語数)を求める問題である。認識精度が最も高くなるような単語列を推定するには、出現確率 $P(W|A)$ が最大となる認識単語列 $W^*$ を求めれば良い。すなわち、

(1)

※(2)式のように書き換える。

★を与えないため、右辺分子を最大化する $W$ を求めれば良い。すなわち、(3)式ようになる。

☆モデルの構築では、あらかじめ多数の話者からの単語、文等の音声データを収集し、統計的手法を利用して認識精度や認識精度と良く関連した指標が向上するように推定する。例えば、音響モデルを構成する隠れマルコフモデルのパラメータを、音響モデルが学習データを出力する尤度が最大となるように、バウム・ウェルチアルゴリズムを用いて推定する。音響モデルの推定方法は、文献1下巻において詳細に記されている。

【0008】同様に、言語モデルの構築では、新聞や会話の書き起こし等のテキストから、言語モデルの構造に従って、それぞれの発話や発話を構成する単語の出現する確率を計算する。例えば、 $n$ グラム言語モデルにおいて、 $n=2$ とおいたとき(バイグラム言語モデルと呼ばれる)、 $P(W)$ は(4)式のように近似される。 $n$ グラム言語モデルのパラメータは、学習用テキストデータ内の隣接する $n$ 単語の頻度から推定される。言語モデルの推定方法は、文献2において詳細に記されている。

ば、文字や形態素、文節等言語的な単位やエントロピー基準に基づいたテキストの分割、これらの組み合わせ等である。

【0010】図12は上記文献1に開示された従来の音声認識装置の構成を示すブロック図である。図12において、101は音声信号100を入力し音声認識して認識結果104を出力する照合手段、102は照合手段1



01が音声認識する際に参照する音響モデル、103は照合手段101が音声認識する際に参照する言語モデルである。

【0011】次に動作について説明する。照合手段101は、ユーザの音声信号100を入力し、音響モデル102及び言語モデル103を参照して、音声認識を実行し認識結果104を出力する。音響モデル102は、入力されたユーザの音声信号100の音声波形を信号処理して得られる音響特徴ベクトルの時系列と、例えば、音素等で表される音声認識装置が扱う最小のシンボル情報との写像関係を表す。言語モデル103は、音響モデル102により写像されるシンボルの組み合わせにより表される単語等より長い認識単位との対応関係と、単語の出現情報を記述する。音響モデル102は、あるシンボルのモデルがベクトル時系列を出力する確率を求めるものであり、すなわち、音声の音響的な観測値系列の確率を求めるものであり、言語モデル103は、ある単語列の出現確率を求めるものである。

【0012】図13は音声信号100を入力して認識結果104を得る従来の音声認識処理の手順を示すフローチャートである。ステップST1301において、入力された音声信号100はA/D変換されてデジタル信号となる。ステップST1302において、デジタル化された音声信号は適当な間隔において信号処理され、音声の性質をよく表す音響特徴ベクトルの時系列へと変換される。

【0013】ステップST1303において、音響特徴ベクトルは音響照合処理により音響モデル102と照合され、それぞれの認識候補について、音響特徴ベクトルの時系列を出力する確率が求められる。ステップST1304において、それぞれの認識候補はさらに言語照合処理によって、言語モデル103と照合され単語列の出力確率が乗じられる。最後にステップST1305において、それぞれの認識候補から最も適切な候補を選択して認識結果104を得る。通常、最も適切な認識結果とは、上記照合によって最も確率が高いとされた認識候補である。

【0014】上記方法により構成された音響モデル及び言語モデルのみでは、十分な性能が達成されない場合で、ユーザがカスタマイズ可能な音声認識装置では、認識されにくい音声や単語をより良く認識させるために、音響モデルをユーザに適応化させたり、認識対象単語をユーザ辞書に追加することによって、認識精度を高めることができる場合がある。

【0015】まず、音響モデルを適応化させる場合について説明する。図14は上記文献1に開示された、音響モデルを音声信号100に適応化させる音響モデル適応化手段を備えた従来の音声認識装置の構成を示すブロック図である。図14において、図12と異なる部分は、入力する音声信号100に対して音響モデルを適応させ

るために、初期音響モデル1003、音響モデル適応化手段1004、適応化済み音響モデル1401を備えることと、音声信号100が音響モデル適応化手段1004及び照合手段101により選択されることである。

【0016】次に図14に示す音声認識装置の動作について説明する。音響モデル適応化手段1004は、実際の認識前に収集した適応化用音声（音声信号100）と初期音響モデル1003から、例えば最大事後確率推定法を用いて、初期音響モデル1003を適応化用音声に適応化させ、適応化済み音響モデル1401を得る。音響モデルの適応化方法については、文献3の7章に示されている。照合手段101は、音声信号100を入力し、適応化済み音響モデル1401及び言語モデル103を参照して、音声認識を行い認識結果104を出力する。

【0017】次に認識対象単語をユーザが辞書に登録する場合について説明する。図15は上記文献1に開示された、ユーザ辞書が追加された従来の音声認識装置の構成を示すブロック図である。図15において、図12と異なる部分はユーザ辞書601を備えることである。

【0018】図16はユーザ辞書601の構成例を示す図である。ユーザ辞書601は認識されない単語、認識されにくい単語を、より良く認識させるためにユーザが登録した単語の集まりで、単語の表記、読みからなる単語の一覧である。

【0019】次に図15に示す音声認識装置の動作について説明する。図17はユーザ辞書601が追加された従来の音声認識処理の手順を示すフローチャートである。図17において、図13と異なる点は、ステップST1704の言語照合処理において、言語モデル103とユーザ辞書601を参照するために、ユーザ辞書601に登録された単語が認識対象単語に加えられることである。

【0020】ユーザ辞書601に登録された単語は、適当な接続確率によって、任意の単語列と接続可能とする。例えば、単語の出現条件を先行する1単語のみにより決定するバイグラム言語モデルにおいて、任意の単語 $w_i$ と $w_j$ には含まれるユーザ辞書登録単語 $w_{user}$ の確率 $P(w_{user}|w_i)$ 、 $P(w_j|w_{user})$ に一定値を与え、言語モデル全体の確率が1になるように確率値を再配分する。この結果、ユーザは登録した単語が含まれた認識結果を得られる。

【0021】しかし、図14及び図15に示す音声認識装置の構成では、認識精度を高めるために、音声信号100と適応化済み音響モデル1401を作成したり、ユーザ辞書601に単語を登録することによって、ユーザが音響モデル102、言語モデル103を自らカスタマイズする必要があり、ユーザに大きな負担を強いることになる。

【0022】また、ユーザ辞書601へ単語を登録した

場合であっても、言語モデル構築時点では、出現していなかった単語や考慮されていなかった用法があるために、それらの単語へ付与した出現確率が不適切である場合がある。さらに、これが原因となり認識精度が低下する可能性がある。

【0023】さらに、上記のようにカスタマイズされた音響モデル102及び言語モデル103は、特定の照合手段101からのみ参照される。このため、それ以外の照合手段101が、このカスタマイズした音響モデル102及び言語モデル103を利用すると、認識精度が低下してしまう。

【0024】このため、漢字かな変換や機械翻訳等の言語処理システムでは、上記に示した問題のうち、ユーザ辞書601の更新に関し、ネットワークを介して辞書を自動的に更新する機能を備えることによって、ユーザの負担を軽減したシステムが、例えば特開平10-260960号公報のように提案されている。しかし、上記公報では、音声のようなパターン認識のためのモデルを扱うことは考慮しておらず、音響モデル102のようなパターン情報に関するモデルに適用することができない。

【0025】さらに、言語モデル103の更新においても、ユーザの登録単語数が増加すると、適切でない短い単語が挿入された湧き出し誤りを生じやすくなることから、認識精度の低下が起きやすくなる。この認識精度の低下は、ユーザ辞書601へ追加登録した単語が、言語モデル103を構築した時点では存在しなかったり、使用環境が変化しているために、例えばユーザ登録単語の出現確率 $P(w_{user}|w_i)$ 、 $P(w_j|w_{user})$ 等に適切でない出現確率が付与されている場合があるからである。この結果、不適切な認識結果を得やすくなり、認識精度が低下する可能性がある。これを防ぐためには、単語の出現確率を適切に設定する必要があるが、ユーザ自身が妥当な値を与えることは一般に困難である。

【0026】

【発明が解決しようとする課題】従来の音声認識装置は以上のように構成されているので、認識精度を高めるために音響モデルや言語モデルをカスタマイズする際に、ユーザに大きな負担をかけるという課題があった。

【0027】また、ユーザごとにカスタマイズされた音声認識装置以外を利用する場合、認識精度が低下するという課題があった。

【0028】さらに、ネットワークを介したカスタマイズによって、音響モデルを自動更新することができないという課題があった。

【0029】さらに、ユーザ辞書への登録が増加すると認識精度が低下しやすいという課題があった。

【0030】この発明は上記のような課題を解決するためになされたもので、ネットワークに接続されたサーバ側で、最新の音響データ又は言語データを取得して、最

新の状態にある音響モデル又は言語モデルを構築したり、ユーザに対応した音響モデル又は言語モデルを構築し、ネットワークを介してユーザ側の音響モデルと言語モデルを更新することで、ユーザに大きな負担をかけることなく認識精度を向上させる音声認識システム、音声認識装置、音響モデル管理サーバ、言語モデル管理サーバ、音声認識方法及び音声認識プログラムを記録したコンピュータ読み取り可能な記録媒体を得ることを目的とする。

10 【0031】また、ネットワークを介して接続することによって、あらゆる音声認識装置において、カスタマイズされた音響モデル又は言語モデルを利用できる音声認識システム、音声認識装置、音響モデル管理サーバ、言語モデル管理サーバ、音声認識方法及び音声認識プログラムを記録したコンピュータ読み取り可能な記録媒体を得ることを目的とする。

【0032】さらに、ユーザから得られる辞書やテキストと半自動的に収集されたテキストを利用することによって、ユーザ辞書が大きくなった場合でも、認識精度が低下しにくい音声認識システム、音声認識装置、音響モデル管理サーバ、言語モデル管理サーバ、音声認識方法及び音声認識プログラムを記録したコンピュータ読み取り可能な記録媒体を得ることを目的とする。

【0033】

【課題を解決するための手段】この発明に係る音声認識システムは、音声信号を入力し、音声の音響的な観測値系列の確率を求める音響モデルを参照して音声認識を行い、認識結果を出力する音声認識装置と、上記音声認識装置とネットワークを介して接続され、更新された音響データを取得して上記音響モデルを構築する音響モデル管理サーバとを備えたものにおいて、上記音響モデル管理サーバが構築した上記音響モデルを上記音声認識装置に送信し、上記音声認識装置が、音声認識の際に参照する音響モデルを、上記音響モデル管理サーバが送信した音響モデルにより更新するものである。

30 【0034】この発明に係る音声認識システムは、音響モデル管理サーバが、音声認識装置が音声認識の際に参照する音響モデルを特定するIDを取得し、取得したIDで指示される特定条件に対応して、更新された音響データを読み出し、上記特定条件に依存した音響モデルを構築して上記音声認識装置に送信するものである。

40 【0035】この発明に係る音声認識システムは、音声信号を入力し、単語列の出現確率を求める言語モデルを参照して音声認識を行い、認識結果を出力する音声認識装置と、上記音声認識装置とネットワークを介して接続され、更新された言語データを取得して上記言語モデルを構築する言語モデル管理サーバとを備えたものにおいて、上記言語モデル管理サーバが構築した上記言語モデルを上記音声認識装置に送信し、上記音声認識装置が、音声認識の際に参照する言語モデルを、上記言語モデル

管理サーバが送信した言語モデルにより更新するものである。

【0036】この発明に係る音声認識システムは、言語モデル管理サーバが、音声認識装置が音声認識の際に参照する言語モデルを特定するIDを取得し、取得したIDで指示される特定条件に対応して、更新された言語データを読み出し、上記特定条件に依存した言語モデルを構築して上記音声認識装置に送信するものである。

【0037】この発明に係る音声認識システムは、音声認識装置が音声認識の際に単語を登録したユーザ辞書を参照し、言語モデル管理サーバが、ネットワークを介して上記ユーザ辞書を読み出し、更新された言語データと、読み出した上記ユーザ辞書とを参照し、上記ユーザ辞書に依存した言語モデルを構築して上記音声認識装置に送信するものである。

【0038】この発明に係る音声認識システムは、言語モデル管理サーバが、音声認識装置のユーザが利用したテキストを取得し、更新された言語データと、取得した上記テキストとを参照し、上記テキストに依存した言語モデルを構築して上記音声認識装置に送信するものである。

【0039】この発明に係る音声認識システムは、音声信号を入力し、音声の音響的な観測値系列の確率を求める音響モデルを参照して音声認識を行い、認識結果を出力する音声認識装置と、上記音声認識装置とネットワークを介して接続され、適応化前の初期音響モデルを有する音響モデル管理サーバとを備えたものにおいて、上記音声認識装置が、上記音響モデルを特定するIDと、入力された音声信号から適応化用の音声データとを取得し、取得したID及び適応化用の音声データを、ネットワークを介して上記音響モデル管理サーバに送信し、上記音響モデル管理サーバが、送信された適応化用の音声データを用いて、上記初期音響モデルを適応化し、適応化済み音響モデルを、送信された上記IDに対応付けて格納すると共に、外部からの音響モデル更新指令を受けて、ネットワークを介して上記音声認識装置から上記音響モデルを特定するIDを受信し、受信したIDに対応する適応化済み音響モデルを、格納している適応化済み音響モデルの中から選択して読み出し、ネットワークを介して上記音声認識装置に送信し、上記音声認識装置が、音声認識の際に参照する音響モデルを、上記音響モデル管理サーバが送信した適応化済み音響モデルにより更新するものである。

【0040】この発明に係る音声認識装置は、音声の音響的な観測値系列の確率を求める音響モデルと、音声信号を入力し上記音響モデルを参照して音声認識を行い、認識結果を出力する照合手段とを備えたものにおいて、ネットワークを介して接続された音響モデル管理サーバから、更新された音響データにより構築された音響モデルを受信し、上記照合手段が音声認識の際に参照する音

響モデルを、受信した音響モデルにより更新する音響モデル更新手段とを備えたものである。

【0041】この発明に係る音声認識装置は、音響モデル更新手段が、ネットワークを介して接続された音響モデル管理サーバから、更新された音響データにより構築された、照合手段が音声認識の際に参照する音響モデルの特定条件に依存した音響モデルを受信し、上記照合手段が音声認識の際に参照する音響モデルを、受信した音響モデルにより更新するものである。

10 【0042】この発明に係る音声認識装置は、単語列の出現確率を求める言語モデルと、音声信号を入力し上記言語モデルを参照して音声認識を行い、認識結果を出力する照合手段とを備えたものにおいて、ネットワークを介して接続された言語モデル管理サーバから、更新された言語データにより構築された言語モデルを受信し、上記照合手段が音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する言語モデル更新手段とを備えたものである。

20 【0043】この発明に係る音声認識装置は、言語モデル更新手段が、ネットワークを介して接続された言語モデル管理サーバから、更新された言語データにより構築された、照合手段が音声認識の際に参照する言語モデルの特定条件に依存した言語モデルを受信し、上記照合手段が音声認識の際に参照する言語モデルを、受信した言語モデルにより更新するものである。

30 【0044】この発明に係る音声認識装置は、照合手段が音声認識の際に参照する単語を登録したユーザ辞書を備え、言語モデル更新手段が、ネットワークを介して接続された言語モデル管理サーバから、更新された言語データにより構築された、上記照合手段が音声認識の際に参照するユーザ辞書に依存した言語モデルを受信し、上記照合手段が音声認識の際に参照する言語モデルを、受信した言語モデルにより更新するものである。

40 【0045】この発明に係る音声認識装置は、言語モデル更新手段が、ネットワークを介して接続された言語モデル管理サーバから、更新された言語データにより構築された、音声認識を行うユーザが利用したテキストに依存した言語モデルを受信し、上記照合手段が音声認識の際に参照する言語モデルを、受信した言語モデルにより更新するものである。

50 【0046】この発明に係る音声認識装置は、音声の音響的な観測値系列の確率を求める音響モデルと、音声信号を入力し上記音響モデルを参照して音声認識を行い、認識結果を出力する照合手段と、上記音響モデルを特定するIDを取得する音響モデルID取得手段と、上記音響モデルID取得手段が取得したIDを読み出し、入力された音声信号から適応化用の音声データを取得し、読み出したID及び取得した適応化用の音声データを、ネットワークを介して接続された音響モデル管理サーバに送信する適応化用音声取得手段と、上記音響モデル管理

サーバから、上記 I D に対応する上記適応化用の音声データにより適応化された適応化済み音響モデルを受信し、上記照合手段が音声認識の際に参照する音響モデルを、受信した適応化済み音響モデルにより更新する音響モデル更新手段とを備えたものである。

【0047】この発明に係る音響モデル管理サーバは、更新された音響データを取得する音響データ取得手段と、外部からの音響モデル更新指令を受けて、上記音響データ取得手段が取得した更新された音響データを読み出し、音声の音響的な観測値系列の確率を求める音響モデルを構築する音響モデル構築手段と、上記音響モデル構築手段により構築された音響モデルを、ネットワークを介して音声認識を行う音声認識装置に送信する音響モデル送信手段とを備えたものである。

【0048】この発明に係る音響モデル管理サーバは、更新された音響データを取得する音響データ取得手段と、外部からの音響モデル更新指令を受けて、ネットワークを介して接続された音声認識装置が音声認識の際に参照する音響モデルを特定する I D を取得する更新音響モデル I D 取得手段と、上記更新音響モデル I D 取得手段が取得した I D で指示される特定条件に対応して、上記音響データ取得手段が取得した更新された音響データを読み出す特定向け音響データ読み出し手段と、上記特定向け音響データ読み出し手段が読み出した更新された音響データを参照し、上記特定条件に依存した音響モデルを構築する特定向け音響モデル構築手段と、上記特定向け音響モデル構築手段が構築した音響モデルを、ネットワークを介して上記音声認識装置に送信する音響モデル送信手段とを備えたものである。

【0049】この発明に係る音響モデル管理サーバは、音声の音響的な観測値系列の確率を求める、適応化前の初期音響モデルと、ネットワークを介して接続された音声認識装置から送信された、適応化用の音声データと、上記音声認識装置が音声認識の際に参照する音響モデルを特定する I D を受信し、上記適応化用の音声データを用いて上記初期音響モデルを適応化し、適応化済み音響モデルを、受信した上記 I D に対応付けて適応化済み音響モデル格納手段に格納する音響モデル適応化手段と、外部からの音響モデル更新指令を受けて、ネットワークを介して上記音声認識装置から上記 I D を受信し、受信した I D に対応する適応化済み音響モデルを、上記適応化済み音響モデル格納手段から選択して読み出す適応化済み音響モデル選択手段と、上記適応化済み音響モデル選択手段が読み出した適応化済み音響モデルを、ネットワークを介して上記音声認識装置に送信する音響モデル送信手段とを備えたものである。

【0050】この発明に係る言語モデル管理サーバは、更新された言語データを取得する言語データ取得手段と、外部からの言語モデル更新指令を受けて、上記言語データ取得手段が取得した更新された言語データを読み

出し、単語列の出現確率を求める言語モデルを構築する言語モデル構築手段と、上記言語モデル構築手段が構築した言語モデルを、ネットワークを介して音声認識を行う音声認識装置に送信する言語モデル送信手段とを備えたものである。

【0051】この発明に係る言語モデル管理サーバは、更新された言語データを取得する言語データ取得手段と、外部からの言語モデル更新指令を受けて、ネットワークを介して接続された音声認識装置が音声認識の際に参照する言語モデルを特定する I D を取得する更新言語モデル I D 取得手段と、上記更新言語モデル I D 取得手段が取得した I D で指示される特定条件に対応して、上記言語データ取得手段が取得した更新された言語データを読み出す特定向け言語データ読み出し手段と、上記特定向け言語データ読み出し手段が読み出した更新された言語データを参照し、上記特定条件に依存した言語モデルを構築する特定向け言語モデル構築手段と、上記特定向け言語モデル構築手段が構築した言語モデルを、ネットワークを介して上記音声認識装置に送信する言語モデル送信手段とを備えたものである。

【0052】この発明に係る言語モデル管理サーバは、更新された言語データを取得する言語データ取得手段と、外部からの言語モデル更新指令を受けて、ネットワークを介して接続された音声認識装置が音声認識の際に参照するユーザ辞書を読み出すユーザ辞書読み出し手段と、上記言語データ取得手段が取得した更新された言語データを読み出し、上記ユーザ辞書読み出し手段が読み出したユーザ辞書に依存した言語モデルを構築するユーザ辞書依存言語モデル構築手段と、上記ユーザ辞書依存言語モデル構築手段が構築した言語モデルを、ネットワークを介して上記音声認識装置に送信する言語モデル送信手段とを備えたものである。

【0053】この発明に係る言語モデル管理サーバは、更新された言語データを取得する言語データ取得手段と、外部からの言語モデル更新指令を受けて、ネットワークを介して接続された音声認識装置のユーザが利用したテキストを取得するユーザ利用テキスト取得手段と、上記言語データ取得手段が取得した更新された言語データを読み出し、上記ユーザ利用テキスト取得手段が取得したテキストに依存した言語モデルを構築するユーザ利用テキスト依存言語モデル構築手段と、上記ユーザ利用テキスト依存言語モデル構築手段が構築した言語モデルを、ネットワークを介して上記音声認識装置に送信する言語モデル送信手段とを備えたものである。

【0054】この発明に係る音声認識方法は、音声信号を入力し、音声の音響的な観測値系列の確率を求める音響モデルを参照して音声認識を行い、認識結果を出力するものにおいて、更新された音響データを取得する第 1 のステップと、音響モデル更新指令を受けて、上記第 1 のステップで取得した更新された音響データを読み出

し、音響モデルを構築する第2のステップと、上記第2のステップで構築した音響モデルを、ネットワークを介して送信する第3のステップと、上記第3のステップで送信した音響モデルを受信し、上記音声認識の際に参照する音響モデルを、受信した音響モデルにより更新する第4のステップとを備えたものである。

【0055】この発明に係る音声認識方法は、音声信号を入力し、単語列の出現確率を求める言語モデルを参照して音声認識を行い、認識結果を出力するものにおいて、更新された言語データを取得する第1のステップと、言語モデル更新指令を受けて、上記第1のステップで取得した更新された言語データを読み出し、言語モデルを構築する第2のステップと、上記第2のステップで構築した言語モデルを、ネットワークを介して送信する第3のステップと、上記第3のステップで送信した言語モデルを受信し、上記音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する第4のステップとを備えたものである。

【0056】この発明に係る音声認識方法は、音声信号を入力し、音声の音響的な観測値系列の確率を求める音響モデルを参照して音声認識を行い、認識結果を出力するものにおいて、更新された音響データを取得する第1のステップと、音響モデル更新指令を受けて、音声認識の際に参照する音響モデルを特定するIDを取得する第2のステップと、上記第2のステップで取得したIDで指示される特定条件に対応して、上記第1のステップで取得した更新された音響データを読み出す第3のステップと、上記第3のステップで読み出した更新された音響データを参照し、上記特定条件に依存した音響モデルを構築する第4のステップと、上記第4のステップで構築した音響モデルを、ネットワークを介して送信する第5のステップと、上記第5のステップで送信した音響モデルを受信し、音声認識の際に参照する音響モデルを、受信した音響モデルにより更新する第6のステップとを備えたものである。

【0057】この発明に係る音声認識方法は、音声信号を入力し、単語列の出現確率を求める言語モデルを参照して音声認識を行い、認識結果を出力するものにおいて、更新された言語データを取得する第1のステップと、言語モデル更新指令を受けて、音声認識の際に参照する言語モデルを特定するIDを取得する第2のステップと、上記第2のステップで取得したIDで指示される特定条件に対応して、上記第1のステップで取得した更新された言語データを読み出す第3のステップと、上記第3のステップで読み出した更新された言語データを参照し、上記特定条件に依存した言語モデルを構築する第4のステップと、上記第4のステップで構築した言語モデルを、ネットワークを介して送信する第5のステップと、上記第5のステップで送信した言語モデルを受信し、音声認識の際に参照する言語モデルを、受信した言

語モデルにより更新する第6のステップとを備えたものである。

【0058】この発明に係る音声認識方法は、音声信号を入力し、単語列の出現確率を求める言語モデルと、単語を登録したユーザ辞書を参照して音声認識を行い、認識結果を出力するものにおいて、更新された言語データを取得する第1のステップと、言語モデル更新指令を受けて、音声認識の際に参照するユーザ辞書を読み出す第2のステップと、上記第1のステップで取得した更新された言語データを読み出し、上記第2のステップで読み出したユーザ辞書に依存した言語モデルを構築する第3のステップと、上記第3のステップで構築した言語モデルを、ネットワークを介して送信する第4のステップと、上記第4のステップで送信した言語モデルを受信し、音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する第5のステップとを備えたものである。

【0059】この発明に係る音声認識方法は、音声信号を入力し、単語列の出現確率を求める言語モデルを参照して音声認識を行い、認識結果を出力するものにおいて、更新された言語データを取得する第1のステップと、言語モデル更新指令を受けて、音声認識を行うユーザが利用したテキストを取得する第2のステップと、上記第1のステップで取得した更新された言語データを読み出し、上記第2のステップで取得したテキストに依存した言語モデルを構築する第3のステップと、上記第3のステップで構築した言語モデルを、ネットワークを介して送信する第4のステップと、上記第4のステップで送信した言語モデルを受信し、音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する第5のステップとを備えたものである。

【0060】この発明に係る音声認識方法は、音声信号を入力し、音声の音響的な観測値系列の確率を求める音響モデルを参照して音声認識を行い、認識結果を出力するものにおいて、上記音響モデルを特定するIDを取得する第1のステップと、上記第1のステップで取得したIDを読み出し、入力された音声信号から適応化用の音声データを取得し、ネットワークを介して、読み出したID及び取得した適応化用の音声データを送信する第2のステップと、上記第2のステップで送信した適応化用の音声データを用いて、適応化前の初期音響モデルを適応化し、適応化済み音響モデルを、上記第2のステップで送信したIDに対応付けて格納する第3のステップと、音響モデル更新指令を受けて、ネットワークを介して上記第1のステップで取得したIDを受信し、受信したIDに対応する適応化済み音響モデルを、上記第3のステップで格納している適応化済み音響モデルの中から選択して読み出す第4のステップと、上記第4のステップで読み出した適応化済み音響モデルを、ネットワークを介して送信する第5のステップと、上記第5のステッ

ブで送信した適応化済み音響モデルを受信し、音声認識の際に参照する音響モデルを、受信した適応化済み音響モデルにより更新する第6のステップとを備えたものである。

【0061】この発明に係る音声認識プログラムを記録したコンピュータ読み取り可能な記録媒体は、音声信号を入力し、音声の音響的な観測値系列の確率を求める音響モデルを参照して音声認識を行い、認識結果を出力する照合機能を実現させるものであって、更新された音響データを取得する音響データ取得機能と、音響モデル更新指令を受けて、上記音響データ取得機能が取得した更新された音響データを読み出し、音響モデルを構築する音響モデル構築機能と、上記音響モデル構築機能が構築した音響モデルを、ネットワークを介して送信する音響モデル送信機能と、上記音響モデル送信機能が送信した音響モデルを受信し、上記照合機能が音声認識の際に参照する音響モデルを、受信した音響モデルにより更新する音響モデル更新機能とを実現させるものである。

【0062】この発明に係る音声認識プログラムを記録したコンピュータ読み取り可能な記録媒体は、音声信号を入力し、単語列の出現確率を求める言語モデルを参照して音声認識を行い、認識結果を出力する照合機能を実現させるものであって、更新された言語データを取得する言語データ取得機能と、言語モデル更新指令を受けて、上記言語データ取得機能が取得した更新された言語データを読み出し、言語モデルを構築する言語モデル構築機能と、上記言語モデル構築機能が構築した言語モデルを、ネットワークを介して送信する言語モデル送信機能と、上記言語モデル送信機能が送信した言語モデルを受信し、上記照合機能が音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する言語モデル更新機能とを実現させるものである。

【0063】この発明に係る音声認識プログラムを記録したコンピュータ読み取り可能な記録媒体は、音声信号を入力し、音声の音響的な観測値系列の確率を求める音響モデルを参照して音声認識を行い、認識結果を出力する照合機能を実現させるものであって、更新された音響データを取得する音響データ取得機能と、音響モデル更新指令を受けて、上記音響モデルを特定するIDを取得する更新音響モデルID取得機能と、上記更新音響モデルID取得機能が取得したIDで指示される特定条件に対応して、上記音響データ取得機能が取得した更新された音響データを読み出す特定向け音響データ読み出し機能と、上記特定向け音響データ読み出し機能が読み出した更新された音響データを参照し、上記特定条件に依存した音響モデルを構築する特定向け音響モデル構築機能と、上記特定向け音響モデル構築機能が構築した音響モデルを、ネットワークを介して送信する音響モデル送信機能と、上記音響モデル送信機能が送信した音響モデルを受信し、上記照合機能が音声認識の際に参照する音響

モデルを、受信した音響モデルにより更新する音響モデル更新機能とを実現させるものである。

【0064】この発明に係る音声認識プログラムを記録したコンピュータ読み取り可能な記録媒体は、音声信号を入力し、単語列の出現確率を求める言語モデルを参照して音声認識を行い、認識結果を出力する照合機能を実現させるものであって、更新された言語データを取得する言語データ取得機能と、言語モデル更新指令を受けて、上記言語モデルを特定するIDを取得する更新言語モデルID取得機能と、上記更新言語モデルID取得機能が取得したIDで指示される特定条件に対応して、上記言語データ取得機能が取得した更新された言語データを読み出す特定向け言語データ読み出し機能と、上記特定向け言語データ読み出し機能が読み出した更新された言語データを参照し、上記特定条件に依存した言語モデルを構築する特定向け言語モデル構築機能と、上記特定向け言語モデル構築機能が構築した言語モデルを、ネットワークを介して送信する言語モデル送信機能と、上記言語モデル送信機能が送信した言語モデルを受信し、上記照合機能が音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する言語モデル更新機能とを実現させるものである。

【0065】この発明に係る音声認識プログラムを記録したコンピュータ読み取り可能な記録媒体は、音声信号を入力し、単語列の出現確率を求める言語モデルと、単語を登録したユーザ辞書を参照して音声認識を行い、認識結果を出力する照合機能を実現させるものであって、更新された言語データを取得する言語データ取得機能と、言語モデル更新指令を受けて、上記ユーザ辞書を読み出すユーザ辞書読み出し機能と、上記言語データ取得機能が取得した更新された言語データを読み出し、上記ユーザ辞書読み出し機能が読み出したユーザ辞書に依存した言語モデルを構築するユーザ辞書依存言語モデル構築機能と、上記ユーザ辞書依存言語モデル構築機能が構築した言語モデルを、ネットワークを介して送信する言語モデル送信機能と、上記言語モデル送信機能が送信した言語モデルを受信し、上記照合機能が音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する言語モデル更新機能とを実現させるものである。

【0066】この発明に係る音声認識プログラムを記録したコンピュータ読み取り可能な記録媒体は、音声信号を入力し、単語列の出現確率を求める言語モデルを参照して音声認識を行い、認識結果を出力する照合機能を実現させるものであって、更新された言語データを取得する言語データ取得機能と、言語モデル更新指令を受けて、音声認識を行うユーザが利用したテキストを取得するユーザ利用テキスト取得機能と、上記言語データ取得機能が取得した更新された言語データを読み出し、上記ユーザ利用テキスト取得機能が取得したテキストに依存した言語モデルを構築するユーザ利用テキスト依存言語

モデル構築機能と、上記ユーザ利用テキスト依存言語モデル構築機能が構築した言語モデルを、ネットワークを介して送信する言語モデル送信機能と、上記言語モデル送信機能が送信した言語モデルを受信し、上記照合機能が音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する言語モデル更新機能とを実現させるものである。

【0067】この発明に係る音声認識プログラムを記録したコンピュータ読み取り可能な記録媒体は、音声信号を入力し、音声の音響的な観測値系列の確率を求める音響モデルを参照して音声認識を行い、認識結果を出力する照合機能を実現させるものであって、上記音響モデルを特定するIDを取得する音響モデルID取得機能と、上記音響モデルID取得機能が取得したIDを読み出し、入力された音声信号から適応化用の音声データを取得し、ネットワークを介して、読み出したID及び取得した適応化用の音声データを送信する適応化用音声取得機能と、上記適応化用音声取得機能が送信した適応化用の音声データを用いて、適応化前の初期音響モデルを適応化し、適応化済み音響モデルを、上記適応化用音声取得機能が送信したIDに対応付けて格納する音響モデル適応化機能と、音響モデル更新指令を受けて、ネットワークを介して上記音響モデルID取得機能が取得したIDを受信し、受信したIDに対応する適応化済み音響モデルを、上記音響モデル適応化機能が格納した適応化済み音響モデルの中から選択して読み出す適応化済み音響モデル選択機能と、上記適応化済み音響モデル選択機能が読み出した適応化済み音響モデルを、ネットワークを介して送信する音響モデル送信機能と、上記音響モデル送信機能が送信した適応化済み音響モデルを受信し、上記照合機能が音声認識の際に参照する音響モデルを、受信した適応化済み音響モデルにより更新する音響モデル更新機能とを実現させるものである。

【0068】

【発明の実施の形態】以下、この発明の実施の一形態を説明する。

実施の形態1. 図1はこの発明の実施の形態1による音声認識システムの構成を示すブロック図である。図において、10は音声認識を行う音声認識装置、20はネットワークに接続されている音響モデル管理サーバ、30はネットワークに接続されている言語モデル管理サーバである。ここで、ネットワークとは、有線あるいは無線によってデジタル信号を伝達可能な通信経路一般を示す。

【0069】音声認識装置10において、100は入力する音声信号、101は音声信号100の音声認識を行う照合手段、102は照合手段101が音声認識の際に参照する音響モデル、103は照合手段101が音声認識の際に参照する言語モデル、104は照合手段101が出力する認識結果、111はネットワークを介して送

信された音響モデルにより音響モデル102を更新する音響モデル更新手段、118はネットワークを介して送信された言語モデルにより言語モデル103を更新する言語モデル更新手段である。

【0070】音響モデル管理サーバ20において、105は外部より与えられる音響モデル更新指令、106は更新された音響データを取得する音響データ取得手段、107は音響データ取得手段106が取得した更新された音響データ、108は音響モデル更新指令105を受けて、更新された音響データ107を読み出し、統計的手法を用いてパラメータ推定を行い音響モデルを構築する音響モデル構築手段、109は構築された音響モデルを格納する音響モデル格納手段、110は音響モデル格納手段109に格納されている音響モデルを、ネットワークを介して音声認識装置10に送信する音響モデル送信手段である。

【0071】言語モデル管理サーバ30において、112は外部より与えられる言語モデル更新指令、113は更新された言語データを取得する言語データ取得手段、114は言語データ取得手段113が取得した更新された言語データ、115は言語モデル更新指令112を受けて、更新された言語データ114を読み出し、統計的手法を用いてパラメータ推定を行い言語モデルを構築する言語モデル構築手段、116は構築された言語モデルを格納する言語モデル格納手段、117は言語モデル格納手段116に格納されている言語モデルを、ネットワークを介して音声認識装置10に送信する言語モデル送信手段である。

【0072】従来技術と異なるこの発明の特徴的な部分は、音響モデル管理サーバ20、言語モデル管理サーバ30において、音響データ取得手段106、言語データ取得手段113により、更新された音響データ107、更新された言語データ114を取得し、音響モデル構築手段108、言語モデル構築手段115により最新の音響モデル、最新の言語モデルを構築し、構築した最新の音響モデル、最新の言語モデルを、ネットワークを介して音声認識装置10に送信し、音声認識装置10において、音響モデル更新手段111、言語モデル更新手段118が、最新の音響モデル、最新の言語モデルにより、照合手段101が参照する音響モデル102、言語モデル103を更新することである。

【0073】次に動作について説明する。音響モデル管理サーバ20において、音響データ取得手段106は、音響モデル更新指令105と同期、あるいは非同期に動作し、常時、更新あるいは配信される音響データを自動、あるいは半自動的にダウンロードし、更新された音響データ107に格納する。取得するこれらの音響データは、例えばインターネット上で更新されたり、マルチメディア放送によって配信される音声データ、あるいは音声データと対応する書き起こしテキストであり、検索



ツールによってインターネット上で検索されたり、マルチメディア放送の番組表を用いて決定されてダウンロードされる。

【0074】更新された音響データ107は、音響データ取得手段106により取得された音響モデル学習用音響データの集積であり、上記の例では、音声データや音声データと対応する書き起こしテキストからなる。

【0075】音響モデル構築手段108は、例えば、一定の時間間隔、音声認識処理が実施された時間間隔、あるいは入力装置から与えられるユーザの指示等、適当なタイミングで与えられる音響モデル更新指令105を受けて、更新された音響データ107を参照し、統計的手法を用いて音響モデルのパラメータ推定を行うことによって、例えば、音声データのみからベクトル量子化アルゴリズムを用いることによって、あるいは音声データと対応する書き起こしテキストからバウム・ウェルチアルゴリズムを用いることによって、学習データを良く表すように音響モデルを構築して音響モデル格納手段109に格納する。

【0076】音響モデル格納手段109は、音響モデル構築手段108により構築された音響モデルを記憶し、読み出し要求に応じて音響モデルを出力する。音響モデル送信手段110は、音響モデル格納手段109から音響モデルを読み出し、音声認識装置10の音響モデル更新手段111に、ネットワークを介して送信する。

【0077】音声認識装置10において、音響モデル更新手段111は、音響モデル管理サーバ20の音響モデル送信手段110からネットワークを介して受け取った音響モデルにより、照合手段101が照合の際に参照する音響モデル102を更新する。

【0078】言語モデル管理サーバ30において、言語データ取得手段113は、言語モデル更新指令112と同期、あるいは非同期に動作し、常時、更新あるいは配信される言語データをダウンロードし、言語モデル構築のために用いる新規の言語データを収集し、更新された言語データ114へ格納する。取得するこれらの言語データは、例えば、定期的に配信される新聞やメールマガジンやインターネット上から検索可能なテキスト、チャット、メール、マニュアル等のテキスト等である。

【0079】更新された言語データ114は、言語データ取得手段113により取得された言語モデル学習用言語データの集積であり、テキストデータや同時に得られるテキスト内容に関するキーワード情報等である。

【0080】言語モデル構築手段115は、例えば、一定の時間間隔、音声認識処理が実施された時間間隔、あるいは入力装置から与えられるユーザの指示等、適当なタイミングで与えられる言語モデル更新指令112を受けて、更新された言語データ114を参照してテキストデータを読み出し、統計的手法を用いて言語モデルのパラメータ推定を行うことによって、例えば単語に分割し

たテキストデータから $n$ グラム統計量を求めることによって、学習用データを良く表すように言語モデルを構築し、言語モデル格納手段116に格納する。

【0081】言語モデル格納手段116は、言語モデル構築手段115により構築された言語モデルを記憶し、読み出し要求に応じて言語モデルを出力する。言語モデル送信手段117は、言語モデル格納手段116から言語モデルを読み出し、ネットワークを介して音声認識装置10の言語モデル更新手段118へ送信する。

【0082】音声認識装置10において、言語モデル更新手段118は、言語モデル管理サーバ30の言語モデル送信手段117からネットワークを介して受け取った言語モデルにより、照合手段101が照合の際に参照する言語モデル103を更新する。

【0083】図2はこの発明の実施の形態1による音響モデル102の更新処理を示すフローチャートである。ステップST201において、音響モデル更新タイミング決定手段(図示していない)は、例えばユーザの指示や音響モデルの最終更新時刻からの時間間隔、ネットワークの利用状況等の監視から適当な更新タイミングを判定し、音響モデル管理サーバ20の音響モデル構築手段108へ音響モデル更新指令105を送信する。音響モデル構築手段108は、音響モデル更新指令105を受けていれば、ステップST202へ進み、音響モデル更新指令105を受けていなければ処理を終了する。

【0084】ステップST202において、音響モデル構築手段108は、学習に用いる更新された音響データ107を読み出す。ステップST203において、音響モデル構築手段108は、更新された音響データ107から、統計的手法を用いて音響モデルのパラメータ推定を行うことによって音響モデルを構築し、構築した音響モデルを音響モデル格納手段109へ格納する。

【0085】ステップST204において、音響モデル送信手段110は、音響モデル格納手段109から音響モデルを読み出して、ネットワークを介して音声認識装置10の音響モデル更新手段111へ音響モデルを送信する。ステップST205において、音響モデル更新手段111は、受け取った音響モデルにより照合手段101が参照する音響モデル102を更新する。

【0086】なお、音響モデル更新指令105により、音響モデル102の更新を要求する際に、同時に音声認識装置10が利用している音響モデル102のバージョンを伝達し、音響モデル送信手段110が、音響モデル格納手段109に格納されている音響モデル全体ではなく、それまで音声認識装置10が利用していた音響モデル102との差分情報のみを送信すれば、送信データを減らすことができ、ネットワークの負荷を軽減することができる。

【0087】また、音響モデル構築手段108が、あらかじめ更新された音響モデルを構築しておき、音響モデ

ル更新指令105の要求にしたがって、音響モデルを音声認識装置10に送信するような形態を取った場合でも同様に動作可能である。

【0088】さらに、音声認識装置10がユーザ辞書を持つ場合であっても、同様に処理可能である。

【0089】さらに、この実施の形態では、音声認識を対象として説明を行ったが、パターンとシンボルの関係を表した確率モデル、シンボルの出現を表した確率モデルからなるパターン認識を対象とするものであれば、同様に適用可能である。

【0090】さらに、更新された音響データ107の格納形式は、音響モデル構築時に利用可能な形式であれば、あらかじめ信号処理や頻度分布を計算してあってもかまわない。

【0091】図3はこの発明の実施の形態1による言語モデル103の更新処理を示すフローチャートである。ステップST301において、言語モデル更新タイミング決定手段(図示していない)は、例えば、ユーザの指示や言語モデルの最終更新時刻からの時間間隔、ネットワークの利用状況等の監視から適当な更新タイミングを判定し、言語モデル管理サーバ30の言語モデル構築手段115に言語モデル更新指令112を送信する。言語モデル構築手段115は、言語モデル更新指令112を受けていればステップST302へ進み、言語モデル更新指令112を受けていなければ処理を終了する。

【0092】ステップST302において、言語モデル構築手段115は、学習に用いる更新された言語データ114を読み出す。ステップST303において、言語モデル構築手段115は、更新された言語データ114から、統計的手法を用いて言語モデルのパラメータ推定を行うことによって言語モデルを構築し、構築した言語モデルを言語モデル格納手段116へ格納する。

【0093】ステップST304において、言語モデル送信手段117は、言語モデル格納手段116から言語モデルを読み出して、ネットワークを介して音声認識装置10の言語モデル更新手段118に送信する。ステップST305において、言語モデル更新手段118は、受け取った言語モデルにより照合手段101が参照する言語モデル103を更新する。

【0094】なお、言語モデル更新指令112により、言語モデルの更新を要求する際に、同時に音声認識装置10が利用している言語モデル103のバージョンを伝達することにより、言語モデル格納手段116に格納されている言語モデル全体ではなく、それまで音声認識装置10が利用していた言語モデル103との差分情報のみを送信すれば、送信データを減らすことができ、ネットワークの負荷を軽減することができる。

【0095】また、言語モデル構築手段115が、あらかじめ更新された言語モデルを構築しておき、言語モデル更新指令112の要求にしたがって、言語モデルを送

信するような形態を取った場合でも、同様に動作可能である。

【0096】さらに、音声認識装置10がユーザ辞書を持つ場合であっても、同様に動作可能である。

【0097】さらに、この実施の形態では、音声認識を対象として説明を行ったが、パターンとシンボルの関係を表した確率モデル、シンボルの出現を表した確率モデルからなるパターン認識を対象とするものであれば、同様に適用可能である。

10 【0098】さらに、更新された言語データ114の格納形式は、言語モデル構築時に利用可能な形式であれば、あらかじめ単語に分割しておいたり、言語モデル構築に使えるように単語や単語連鎖、同時に出現する単語の組み合わせ等について頻度あるいは確率計算してあってもかまわない。

【0099】この実施の形態1の図1では、音響モデル更新指令105に従って音響モデルを構築する場合について示しているが、音響データ取得手段106が音響データを取得した際に、音響モデル構築手段108が音響モデルを更新して音響モデル格納手段109に格納し、音響モデル送信手段110が音響モデル更新指令105を受けて音響モデルを読み出すようにしても良い。言語モデルの更新についても同様である。

【0100】また、図1では音響モデル更新手段111及び言語モデル更新手段118を備える場合を示したが、音響モデル更新手段111又は言語モデル更新手段118のどちらか一方のみを備える場合であってもかまわない。

【0101】なお、実施の形態1における音声認識システムを音声認識プログラムとして記録媒体に記録することもできる。この場合には、音響モデル管理サーバ20において、音響データ取得手段106と同様の処理を行う音響データ取得機能と、音響モデル構築手段108と同様の処理を行う音響モデル構築機能と、音響モデル格納手段109と同様の処理を行う音響モデル格納機能と、音響モデル送信手段110と同様の処理を行う音響モデル送信機能から構成されるソフトウェアと、言語モデル管理サーバ30において、言語データ取得手段113と同様の処理を行う言語データ取得機能と、言語モデル構築手段115と同様の処理を行う言語モデル構築機能と、言語モデル格納手段116と同様の処理を行う言語モデル格納機能と、言語モデル送信手段117と同様の処理を行う言語モデル送信機能から構成されるソフトウェアと、音声認識装置10において、音響モデル更新手段111と同様の処理を行う音響モデル更新機能と、言語モデル更新手段118と同様の処理を行う言語モデル更新機能と、照合手段101と同様の処理を行う照合機能から構成されるソフトウェアで音声認識プログラムとなる。

50 【0102】記録媒体に記録する音声認識プログラム

は、音声認識装置10のソフトウェアと、音響モデル管理サーバ20のソフトウェアを、別々の記録媒体に記録しても良いし、1つの記録媒体に記録して、音声認識装置10又は音響モデル管理サーバ20から、それぞれ音響モデル管理サーバ20又は音声認識装置10に送信しても良い。また、これは言語モデルを対象とした場合でも同様である。

【0103】以上のように、この実施の形態1によれば、ネットワークに接続された音響モデル管理サーバ20又は言語モデル管理サーバ30で、更新された音響データ107又は更新された言語データ114を取得し、最新の状態にある音響モデル又は言語モデルを構築し、ネットワークを介してユーザ側の音声認識装置10の音響モデル102又は言語モデル103を更新することで、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果が得られる。

【0104】実施の形態2. 図4はこの発明の実施の形態2による音声認識システムの構成を示すブロック図である。図4の言語モデル管理サーバ30において、401は更新する音声認識装置10の言語モデル103を特定するIDを取得する更新言語モデルID取得手段、402は、更新言語モデルID取得手段401が取得したIDで指示される特定条件に対応して、更新された言語データ114を読み出す特定向け言語データ読み出し手段、403は、特定向け言語データ読み出し手段402により読み出された更新された言語データ114を参照し、特定条件に対応した言語モデルを構築する特定向け言語モデル構築手段である。既に説明した各手段及び各モデルについては、同一の符号を付し説明を省略する。

【0105】従来技術と異なるこの実施の形態に特徴的な部分は、更新言語モデルID取得手段401と、特定向け言語データ読み出し手段402と、特定向け言語モデル構築手段403とを備え、ネットワークを介して言語モデルIDにより特定された言語モデルを提供することである。ここで、特定向け言語モデルとは、特定のユーザやグループ、応用アプリケーション等に言語モデルを特化させることによって、より高い性能が得られるように学習した言語モデルである。

【0106】次に動作について説明する。言語管理サーバ30において、更新言語モデルID取得手段401は、更新された複数の言語モデル103から特定の言語モデル103を選択するために用いられるIDを取得する。このIDは、例えば利用者のユーザID、音声信号100の対象となるタスクを表すID等であり、更新する特定向けの言語モデル103を一意に定めることができるものである。

【0107】特定向け言語データ読み出し手段402は、更新言語モデルID取得手段401が取得した更新言語モデルIDを受け取り、更新された言語データ11

4を文や独立したテキストの単位で読み出して、例えば言語データに付与されるテキスト内容に関するキーワードや言語データに含まれるキーワードから判定し、言語モデルIDで特定される対象であるかどうか識別するフラグを付与し、特定向け言語モデル構築手段403へ送る。

【0108】特定向け言語モデル構築手段403は、更新された言語データ114から特定対象について認識精度が高くなるように学習した言語モデルを構築し、言語モデル格納手段116に格納する。

【0109】このために、まず、学習用言語データを文、あるいは複数の文を単位として含まれるキーワード等から、特定される言語モデル103を判定しておき、これに従って、例えば「対話音声認識のための事前タスク適応の検討」、伊藤彰則、好田正紀、電子情報通信学会技術研究報告(SP96-81)、1996年(以下、文献4とする)で検討されているように、関連深いテキストデータにより大きな重みを付与することによって、特定条件に対応した言語モデルを構築できる。

【0110】例えば、スポーツに関するトピックに特化した特定向け言語モデルを構築するには、言語モデル学習時に、特定向け言語データ読み出し手段402から得たフラグを参照し、スポーツに関するトピックのテキストデータであれば、実際の頻度を $\alpha$ 倍して数え、それ以外の記事であれば、そのまま頻度で両者を加えて確率モデルを推定する。ここで、 $\alpha$ は音声認識の対象とする特定向けテキストのうち、学習に用いないデータに対する言語モデルのエントロピーが最小となるように定める。

【0111】図5はこの発明の実施の形態2による言語モデル103の更新処理を示すフローチャートである。ステップST501において、言語モデル更新タイミング決定手段(図示していない)は、例えば、ユーザの指示や最終更新時刻からの時間間隔、ネットワークの利用状況等の監視から適当な更新タイミングを判定し、言語モデル管理サーバ30の更新言語モデルID取得手段401へ言語モデル更新指令112を送信する。更新言語モデルID取得手段401は、言語モデル更新指令112を受けていれば、ステップST502へ進み、言語モデル更新指令112を受けていなければ、処理を終了する。

【0112】ステップST502において、更新タイミングであれば、更新言語モデルID取得手段401は、使用しているユーザ・グループを特定する手段、タスクを特定する手段等により、更新先の言語モデル103のIDを取得し、特定向け言語データ読み出し手段402へ送る。ステップST503において、特定向け言語データ読み出し手段402は、特定向け言語モデルIDに従い、更新された言語データ114を文や独立したテキストの単位で読み出して、言語モデルIDで特定される対象であるかどうか判別するフラグを付与して、更新さ

れた言語データ114を読み出す。

【0113】ステップST504において、特定向け言語モデル構築手段403は、学習アルゴリズムに従い、更新された言語データ114から言語モデルを推定し、推定した言語モデルを言語モデル格納手段116へ格納する。ステップST505において、言語モデル送信手段117は、読み出した言語モデルを、ネットワークを介して音声認識装置10の言語モデル更新手段118に送信する。ステップST506において、言語モデル更新手段118は、受け取った言語モデルにより照合手段101が参照する言語モデル103を更新する。

【0114】なお、言語モデルIDに対して、適合する言語モデルを出力することができれば、特定向け言語モデルをあらかじめ構築しておく必要はない。例えば、言語モデルIDに依存して学習用言語データのみ作成しておき、要求に応じて言語モデルを構築しても良い。

【0115】なお、この実施の形態では、文献4にしたがった特定向け言語モデル構成法を例としたが、言語モデルを決定できるIDを用いて、複数の言語モデル103から選択する方法であれば同様に適用可能である。

【0116】また、言語モデル更新指令112により、言語モデル更新要求の際に、同時に音声認識装置10が利用している言語モデル103のバージョンを伝達することにより、構築した特定向け言語モデル全体ではなく、現行の言語モデル103からの差分情報のみを送信し、ネットワークの負荷を軽減することができる。

【0117】さらに、音声認識装置にユーザ辞書601がある場合であっても、同様に処理可能である。

【0118】さらに、この説明では音声認識を対象として説明を行ったが、パターンとシンボルの関係を表した確率モデル、シンボルの出現を表した確率モデルからなるパターン認識を対象とするものであれば同様に適用可能である。

【0119】また、この実施の形態では、特定向け言語モデルを構築し、構築した特定向け言語モデルにより言語モデル103を更新しているが、特定向け音響モデルを構築し、構築した特定向け音響モデルにより音響モデル102を更新することもできる。その場合には、図4において、言語モデル管理サーバ30の代わりに音響モデル管理サーバ、言語データ取得手段113の代わりに音響データ取得手段、更新された言語データ114の代わりに更新された音響データ、更新言語モデルID取得手段401の代わりに更新音響モデルID取得手段、特定向け言語データ読み出し手段402の代わりに特定向け音響モデル読み出し手段、特定向け言語モデル構築手段403の代わりに特定向け音響モデル構築手段、言語モデル格納手段116の代わりに音響モデル格納手段、言語モデル送信手段117の代わりに音響モデル送信手段を備え、音声認識装置10において、言語モデル更新手段118の代わりに音響モデル更新手段を備え、音響

モデル更新手段が音響モデル102を更新するようにすれば良い。

【0120】さらに、実施の形態2における音声認識システムを音声認識プログラムとして記録媒体に記録することもできる。この場合には、言語モデル管理サーバ30において、言語データ取得手段113と同様の処理を行う言語データ取得機能と、更新言語モデルID取得手段401と同様の処理を行う更新言語モデルID取得機能と、特定向け言語データ読み出し手段402と同様の処理を行う特定向け言語データ読み出し機能と、特定向け言語モデル構築手段403と同様の処理を行う特定向け言語モデル構築機能と、言語モデル格納手段116と同様の処理を行う言語モデル格納機能と、言語モデル送信手段117と同様の処理を行う言語モデル送信機能から構成されるソフトウェアと、音声認識装置10において、言語モデル更新手段118と同様の処理を行う言語モデル更新機能と、照合手段101と同様の処理を行う照合機能から構成されるソフトウェアで音声認識プログラムとなる。これは音響モデルを対象とした場合でも同様である。

【0121】記録媒体に記録する音声認識プログラムは、音声認識装置10のソフトウェアと、言語モデル管理サーバ30のソフトウェアを、別々の記録媒体に記録しても良いし、1つの記録媒体に記録して、音声認識装置10又は言語モデル管理サーバ30から、それぞれ言語モデル管理サーバ30又は音声認識装置10に送信しても良い。また、これは音響モデルを対象とした場合でも同様である。

【0122】以上のように、この実施の形態2によれば、ネットワークに接続された言語モデル管理サーバ30で、更新された言語データ114を取得し、しかも、ユーザの音声認識装置10の言語モデル103のIDを取得することで、最新の状態にあり、しかもユーザに対応した特定向けの言語モデルを構築し、ネットワークを介してユーザの音声認識装置10の言語モデル103を更新することで、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果が得られる。

【0123】また、この実施の形態2によれば、特定向けにカスタマイズされた言語モデルを、ネットワークを介して更新することにより、ユーザが複数の異なる照合手段101を利用する場合でも、全ての照合手段101の利用時に適切な言語モデル103を利用し、高い認識精度を得ることができるという効果が得られる。

【0124】実施の形態3。図6はこの発明の実施の形態3による音声認識システムの構成を示すブロック図である。図6の音声認識装置10において、601は照合手段101が照合の際に参照する単語を登録したユーザ辞書であり、言語モデル管理サーバ30において、602は言語モデル更新指令112を受けて、照合手段10

1が参照するユーザ辞書601を、ネットワークを介して読み出すユーザ辞書読み出し手段603は、更新された言語データ114とユーザ辞書読み出し手段602が読み出したユーザ辞書601を参照し、ユーザ辞書601に依存した言語モデルを構築するユーザ辞書依存言語モデル構築手段である。既に説明した各手段及び各モデルについては、同一の符号を付し説明を省略する。

【0125】従来技術と異なるこの実施の形態に特徴的な部分は、ユーザ辞書読み出し手段602とユーザ辞書依存言語モデル構築手段603を備えたことである。

【0126】次に動作について説明する。言語モデル管理サーバ30のユーザ辞書読み出し手段602は、言語モデル更新指令112を受けて、音声認識装置10の照合手段101が参照するユーザ辞書601を、ネットワークを介して読み出す。ユーザ辞書依存言語モデル構築手段603は、ユーザ辞書601に登録された単語と更新された言語データ114を用いて、最新の状態に更新されており、かつユーザにカスタマイズされた言語モデルを構築する。

【0127】ユーザ辞書601に依存した言語モデルの構築は、例えば、更新された言語データ114の中から、ユーザ辞書601に存在する単語が含まれるテキストを抜き出し、これを特定向けテキストであると見なし、実施の形態2で参照した文献4記載の方法を実施することにより行われる。これによって、ユーザ辞書601記載の単語のうち、元の言語モデルでは登録されておらず、適切な統計量が付与されていなかった単語で、更新されたテキストにおいて出現した単語に妥当な統計量を付与することが可能となり、認識精度が向上することを期待できる。

【0128】図7はこの発明の実施の形態3による言語モデル103の更新処理を示すフローチャートである。ステップST701において、言語モデル更新タイミング決定手段（図示していない）は、例えば、ユーザの指示や最終更新時刻からの時間間隔、ネットワークの利用状況等の監視から適当な更新タイミングを判定し、言語モデル管理サーバ30のユーザ辞書読み出し手段602へ言語モデル更新指令112を送信する。ユーザ辞書読み出し手段602は、言語モデル更新指令112を受けていれば、ステップST702へ進み、言語モデル更新指令112を受けていなければ、処理を終了する。

【0129】ステップST702において、ユーザ辞書読み出し手段602は、照合手段101が参照するユーザ辞書601を、ネットワークを介して読み出す。ステップST703において、ユーザ辞書依存言語モデル構築手段603は、さらに更新された言語データ114を読み出す。ステップST704において、ユーザ辞書依存言語モデル構築手段603は、ユーザ辞書601及び更新された言語データ114からユーザ辞書601に依存した言語モデルを構築し、言語モデル格納手段116

に格納する。

【0130】ステップST705において、言語モデル送信手段117は、言語モデル格納手段116から読み出したユーザ辞書601に依存した言語モデルを、ネットワークを介して音声認識装置10の言語モデル更新手段118に送信する。ステップST706において、言語モデル更新手段118は、受け取った言語モデルにより照合手段101が参照する言語モデル103を更新する。

【0131】また、この実施の形態では、音声認識を対象として説明を行ったが、パターンとシンボルの関係を表した確率モデル、シンボルの出現を表した確率モデルからなるパターン認識を対象とするものであれば、同様に適用可能である。

【0132】さらに、実施の形態3における音声認識システムを音声認識プログラムとして記録媒体に記録することもできる。この場合には、言語モデル管理サーバ30において、言語データ取得手段113と同様の処理を行う言語データ取得機能と、ユーザ辞書読み出し手段602と同様の処理を行うユーザ辞書読み出し機能と、ユーザ辞書依存言語モデル構築手段603と同様の処理を行うユーザ辞書依存言語モデル構築機能と、言語モデル格納手段116と同様の処理を行う言語モデル格納機能と、言語モデル送信手段117と同様の処理を行う言語モデル送信機能から構成されるソフトウェアと、音声認識装置10において、言語モデル更新手段118と同様の処理を行う言語モデル更新機能と、照合手段101と同様の処理を行う照合機能から構成されるソフトウェアで音声認識プログラムとなる。

【0133】記録媒体に記録する音声認識プログラムは、音声認識装置10のソフトウェアと、言語モデル管理サーバ30のソフトウェアを、別々の記録媒体に記録しても良いし、1つの記録媒体に記録して、音声認識装置10又は言語モデル管理サーバ30から、それぞれ言語モデル管理サーバ30又は音声認識装置10に送信しても良い。

【0134】以上のように、この実施の形態3によれば、ネットワークに接続された言語モデル管理サーバ30で、更新された言語データ114を取得し、しかも、ユーザの音声認識装置10のユーザ辞書601を読み出すことで、最新の状態にあり、しかもユーザ辞書601に登録した単語について、より詳細に反映させた言語モデルを構築し、ネットワークを介してユーザの音声認識装置10の言語モデル103を更新することで、ユーザ辞書が大きくなった場合でも、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果が得られる。

【0135】実施の形態4。図8はこの発明の実施の形態4による音声認識システムの構成を示すブロック図である。図8の言語モデル管理サーバ30において、80

1は、言語モデル更新指令112を受けて、ユーザが利用したテキストを取得するユーザ利用テキスト取得手段、802はユーザ利用テキスト取得手段801が取得したテキストを格納するユーザ利用テキスト格納手段、803は、更新された言語データ114と、ユーザ利用テキスト格納手段802に格納されているテキストを参照し、テキストに依存した言語モデルを構築するユーザ利用テキスト依存言語モデル構築手段である。既に説明した各手段及び各モデルについては、同一の符号を付し説明を省略する。

【0136】従来技術と異なるこの実施の形態に特徴的な部分は、ユーザ利用テキスト取得手段801、ユーザ利用テキスト格納手段802及びユーザ利用テキスト依存言語モデル構築手段803を備え、ユーザが利用したテキストと最新状態に更新された言語データ114を参照し、ユーザが利用したテキストに合わせて言語モデルを構築することである。

【0137】次に動作について説明する。言語モデル管理サーバ30のユーザ利用テキスト取得手段801は、言語モデル更新指令112を受けて、例えば、ユーザがあらかじめ指定したファイル、ディレクトリを走査することにより、ユーザが参照、あるいは記述したテキストファイルを読み出す。ユーザ利用テキスト格納手段802は、ユーザ利用テキスト取得手段801によって収集されたテキストを格納する。

【0138】ユーザ利用テキスト依存言語モデル構築手段803は、ユーザ利用テキスト及び更新された言語データ114を参照し、認識精度が高くなるように言語モデルを構築する。ユーザ利用テキストを用いた言語モデルの構築では、例えば、ユーザ利用テキストを特定向けテキストであると見なし、実施の形態2で参照した文献4記載の方法を実施することにより、ユーザ利用テキスト依存の言語モデルを構築する。このようにして構築された言語モデルは、ユーザが参照あるいは既出したテキストの性質を反映させているため、ユーザが発声する確率の高い言語的性質を含み、より精度の高い認識結果を得ることができる。

【0139】図9はこの発明の実施の形態4による言語モデル103の更新処理を示すフローチャートである。ステップST901において、言語モデル更新タイミング決定手段(図示していない)は、ユーザの指示や最終更新時刻からの時間間隔、ネットワークの利用状況等のモニタから適当なタイミングを判定し、言語モデル管理サーバ30のユーザ利用テキスト取得手段801へ言語モデル更新指令112を送信する。ユーザ利用テキスト取得手段801は、言語モデル更新指令112を受けていれば、ステップST902へ進み、言語モデル更新指令112を受けていなければ、処理を終了する。

【0140】ステップST902において、ユーザ利用テキスト取得手段801は、ユーザ利用テキストを読み

出しユーザ利用テキスト格納手段802へ格納する。ステップST903において、ユーザ利用テキスト依存言語モデル構築手段803は、ユーザ利用テキストと更新された言語データ114を読み出す。ステップST904において、ユーザ利用テキスト依存言語モデル構築手段803は、ユーザ利用テキスト及び更新された言語データ114からユーザ利用テキスト依存言語モデルを構築し言語モデル格納手段116に格納する。

【0141】ステップST905において、言語モデル送信手段117は、言語モデル格納手段116から読み出したユーザ利用テキスト依存言語モデルを、ネットワークを介して音声認識装置10の言語モデル更新手段118に送信する。ステップST906において、言語モデル更新手段118は、受け取った言語モデルにより照合手段101が参照する言語モデル103を更新する。

【0142】この実施の形態では、ユーザ利用テキストの入手を、特定ディレクトリやファイルの検索によるとしたが、ユーザ利用テキスト取得手段801は、テキストが収集できるのであれば、テキストをファイルやディレクトリから取り出すのではなく、音声認識やキーボード、ペン、OCR等のユーザ入力、あるいはブラウザ等によってユーザが閲覧したテキスト等を利用してもかまわない。

【0143】また、ユーザ利用テキスト格納手段802は、ユーザ利用テキスト取得手段801によって収集されたテキストを格納するとしたが、ユーザ利用テキスト依存言語モデル構築手段803の基準に従い、テキストを適当な手段によって単語に分割したり、モデル構築時に参照する単語、単語連鎖、同時に出現する単語の組み合わせ等に関する頻度として格納しても同様である。

【0144】さらに、この実施の形態では、音声認識を対象として説明を行ったが、パターンとシンボルの関係を表した確率モデル、シンボルの出現を表した確率モデルからなるパターン認識を対象とするものであれば同様に適用可能である。

【0145】さらに、実施の形態4における音声認識システムを、音声認識プログラムとして記録媒体に記録することもできる。この場合には、言語モデル管理サーバ30において、言語データ取得手段と同様の処理を行う言語データ取得機能と、ユーザ利用テキスト取得手段801と同様の処理を行うユーザ利用テキスト取得機能と、ユーザ利用テキスト格納手段802と同様の処理を行うユーザ利用テキスト格納機能と、ユーザ利用テキスト依存言語モデル構築手段803と同様の処理を行うユーザ利用テキスト依存言語モデル構築機能と、言語モデル格納手段116と同様の処理を行う言語モデル格納と、言語モデル送信手段117と同様の処理を行う言語モデル送信機能から構成されるソフトウェアと、音声認識装置10において、言語モデル更新手段118と同様の処理を行う言語モデル更新機能と、照合手段101と



同様の処理を行う照合機能から構成されるソフトウェアで音声認識プログラムとなる。

【0146】記録媒体に記録する音声認識プログラムは、音声認識装置10のソフトウェアと、言語モデル管理サーバ30のソフトウェアを、別々の記録媒体に記録しても良いし、1つの記録媒体に記録して、音声認識装置10又は言語モデル管理サーバ30から、それぞれ言語モデル管理サーバ30又は音声認識装置10に送信しても良い。

【0147】以上のように、この実施の形態4によれば、ネットワークに接続された言語モデル管理サーバ30で、更新された言語データ114を取得し、しかも、ユーザが利用したテキストを取得することで、最新の状態にあり、しかもユーザが利用したテキストに依存した言語モデルを構築し、ネットワークを介してユーザの音声認識装置10の言語モデル103を更新することで、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果が得られる。

【0148】実施の形態5. 図10はこの発明の実施の形態5による音声認識システムの構成を示すブロック図である。図10の音声認識装置10において、1001は照合手段101が照合の際に参照する音響モデルを識別するIDを取得する音響モデルID取得手段、1002は、音響モデルID取得手段1001が取得したIDを読み込み、入力された音声信号から適応化用の音声データを取得し、読み込んだIDと取得した適応化用の音声データを、ネットワークを介して音響モデル管理サーバ20に送信する適応化用音声取得手段である。

【0149】図10の音響モデル管理サーバ20において、1003は適応化前の初期音響モデル、1004は、音声認識装置10の適応化用音声取得手段1002から送信された適応化用の音声データを用いて、適応化前の初期音響モデル1003を適応化し、適応化済み音響モデルを、適応化用音声取得手段1002から送信されたIDに対応付けて適応化済み音響モデル格納手段1005に格納する音響モデル適応化手段、1006は、音響モデル更新指令105を受けて、ネットワークを介して音声認識装置10の音響モデルID取得手段1001が取得したIDを受信し、受信したIDに対応する適応化済み音響モデルを、適応化済み音響モデル格納手段1005から選択して読み出す適応化済み音響モデル選択手段である。既に説明した各手段及び各モデルについては、同一の符号を付し説明を省略する。

【0150】従来技術と異なるこの実施の形態に特徴的な部分は、音声認識装置10において、音響モデルID取得手段1001、適応化用音声取得手段1002を備え、音響モデル管理サーバ20において、初期音響モデル1003、音響モデル適応化手段1004、適応化済み音響モデル格納手段1005及び適応化済み音響モデル選択手段1006を備えたことである。

【0151】この実施の形態では、音声認識装置10において、適応化対象となる音響モデルIDと適応化用の音声データを取得し、音響モデル管理サーバ20において、音響モデルIDに依存した適応化を行った音響モデルを構築して、音声認識装置10に送信し、音声認識装置10において、音響モデル102を更新することにより、ユーザは任意の照合手段101の利用に際して、適応化した音響モデルを参照可能であるため、より高い認識精度を得ることができる。

10 【0152】次に動作について説明する。音声認識装置10において、音響モデルID取得手段1001は、適応化対象となる音響モデルを決定するものであり、例えば、音声認識装置10を利用するユーザのユーザIDである。適応化用音声取得手段1002は、照合手段101の利用前にあらかじめ適応化用として音声信号100による適応化用の音声データを取得し、音響モデルID取得手段1001から読み出す音響モデルIDと、取得した適応化用の音声データを、ネットワークを介して接続される音響モデル管理サーバ20の音響モデル適応化手段1004へ送信する。

【0153】音響モデル管理サーバ20において、初期音響モデル1003は適応化を行う前の音響モデルである。音響モデル適応化手段1004は、ネットワークを介して受け取った適応化用の音声データと初期音響モデル1003を用いて、適応化した音響モデルを構築し、適応化済みの音響モデルとネットワークを介して受け取った音響モデルIDを、適応化済み音響モデル格納手段1005へ格納する。音響モデルの適応化には、例えば最大事後確率推定法を用いる。

30 【0154】適応化済み音響モデル格納手段1005は、音響モデルIDと音響モデル適応化手段1004により適応化された音響モデルを格納し、適応化済み音響モデル選択手段の要求に従い、指定の音響モデルIDを持つ音響モデルを出力する。適応化済み音響モデル選択手段1006は、音響モデル更新指令105を受けて、音響モデルID取得手段1001から音響モデルIDを取得し、対応する適応化済み音響モデルを、適応化済み音響モデル格納手段1005から選択して読み出す。

40 【0155】図11はこの発明の実施の形態5による音響モデル102の更新処理を示すフローチャートである。処理はステップST1101からST1107までの音響モデルの適応化段階と、ステップST1108からST1112までの音響モデルの更新段階に分けられる。適応化段階では、入力された適応化用の音声データを用いて、音響モデルIDに依存した音響モデルを構築する。

50 【0156】ステップST1101において、音声認識装置10の音響モデルID取得手段1001は、適応化対象となる音響モデルを識別する、例えばユーザ名等の識別情報を取得する。ステップST1102において、

適応化用音声取得手段1002は、音響モデルIDを読み出し、音声信号100から入力される適応化用の音声データを取得する。ステップST1103において、適応化用音声取得手段1002は、ネットワークを介して音響モデル管理サーバ20に音響モデルの適応化要求を送信し、同時に音響モデルIDと適応化用の音声データを音響モデル適応化手段1004へ送信する。

【0157】ステップST1104において、音響モデル適応化手段1004は、ネットワークを介して音響モデルの適応化要求を受信し、音響モデルIDと適応化用の音声データを読み出す。ステップST1105において、音響モデル適応化手段1004は初期音響モデル1003を読み出す。ステップST1106において、音響モデル適応化手段1004は、ネットワークを介して受け取った適応化用の音声データを用いて初期音響モデル1003を適応化する。ステップST1107において、音響モデル適応化手段1004は、適応化された音響モデルを音響モデルIDによって区別できるように、適応化済み音響モデル格納手段1005に格納する。

【0158】ステップST1108において、音響モデル更新タイミング決定手段（図示していない）は、ユーザの指示や最終更新時刻からの時間間隔、ネットワークの利用状況等の監視から適当な更新タイミングを判定し、音響モデル更新指令105を適応化済み音響モデル選択手段1006へ送信する。適応化済み音響モデル選択手段1006は、音響モデル更新指令105を受けていれば、ステップST1109へ進み、音響モデル更新指令105を受けていなければ処理を終了する。ステップST1109において、適応化済み音響モデル選択手段1006は、音声認識装置10の音響モデルID取得手段1001から、ネットワークを介して適応化対象の音響モデルIDを読み出す。

【0159】ステップST1110において、適応化済み音響モデル選択手段1006は、音響モデルIDで指定された音響モデルを、適応化済み音響モデル格納手段1005から選択して読み出す。ステップST1111において、音響モデル送信手段110は、読み出された適応化済み音響モデルを、ネットワークを介して音声認識装置10の音響モデル更新手段111に送信する。ステップST1112において、音響モデル更新手段111は、受け取った適応化済み音響モデルにより照合手段101が参照する音響モデル102を更新する。

【0160】なお、音声認識装置10において、音響モデルID取得手段1001は、適応化対象となる音響モデルを決定するものであれば、回線伝達特性、背景雑音特性、残響音特性等を決定するものであってもかまわない。

【0161】また、音声認識装置10において、適応化用音声取得手段1002は、照合手段101の利用前にあらかじめ適応化用としてユーザの音声データを取得し

ているが、照合手段101が照合時に音声データを取得してして適応化用音声取得手段1002に入力することで、次の照合の際に参照する音響モデルの適応化を行うことも可能である。

【0162】さらに、音声認識装置10において、適応化用の音声のデータの取得と音響モデルIDの取得の手順は逆であってもかまわない。

【0163】さらに、音声認識装置10において、適応化用音声取得手段1002の音声データの格納形態は音声波形としたが、音声データを信号処理した音響特徴ベクトルの時系列、音響特徴ベクトルとベクトル量子化コードブックを参照して得られるコードブック符号列、それらを統計処理して得られる頻度分布、頻度分布から得られる確率分布等、音響モデルの学習に利用できる形態であれば、どのようなものでもかまわない。

【0164】さらに、音響モデル管理サーバ20において、音響モデル適応化手段1004が受信した適応化用の音声データを格納する手段を追加し、格納した多くの適応化用の音声データにより初期モデル1003を更新することにより、学習の精度が高くなり、より高精度な認識を行う適応化済み音響モデルを構築することができる。

【0165】さらに、この実施の形態では音声認識を対象としたが、パターンとシンボルの関係を表した確率モデル、シンボルの出現を表した確率モデルからなるパターン認識を対象とするものであれば同様に適用可能である。

【0166】さらに、実施の形態5における音声認識システムを音声認識プログラムとして記録媒体に記録することもできる。この場合には、音声認識装置10において、音響モデルID取得手段1001と同様の処理を行う音響モデルID取得機能と、適応化用音声取得手段1002と同様の処理を行う適応化用音声取得機能と、音響モデル更新手段111と同様の処理を行う音響モデル更新機能と、照合手段101と同様の処理を行う照合機能から構成されるソフトウェアと、音響モデル管理サーバ20において、音響モデル適応化手段1004と同様の処理を行う音響モデル適応化機能と、適応化済み音響モデル格納手段1005と同様の処理を行う適応化済み音響モデル格納機能と、適応化済み音響モデル選択手段1006と同様の処理を行う適応化済み音響モデル選択機能と、音響モデル送信手段110と同様の処理を行う音響モデル送信機能から構成されるソフトウェアで音声認識プログラムとなる。

【0167】記録媒体に記録する音声認識プログラムは、音声認識装置10のソフトウェアと、音響モデル管理サーバ20のソフトウェアを、別々の記録媒体に記録しても良いし、1つの記録媒体に記録して、音声認識装置10又は音響モデル管理サーバ20から、それぞれ音響モデル管理サーバ20又は音声認識装置10に送信し

でも良い。

【0168】以上のように、この実施の形態5によれば、ネットワークに接続された音響言語モデル管理サーバ20で、ユーザの音声信号100による適応化用の音声データにより適応化済みの音響モデルを構築し、ネットワークを介してユーザの音声認識装置10の音響モデル102を更新することで、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果が得られる。

【0169】また、この実施の形態5によれば、ユーザに適応化した適応化済み音響モデルにより、ネットワークを介して音響モデル102を更新することで、ユーザが複数の異なる照合手段101を利用する場合でも、全ての照合手段101の利用時に適切な音響モデル102を利用し、高い認識精度を得ることができるという効果が得られる。

【0170】

【発明の効果】以上のように、この発明によれば、音響モデル管理サーバが、更新された音響データを取得して構築した音響モデルを、ネットワークを介して音声認識装置に送信し、音声認識装置が、音声認識の際に参照する音響モデルを、音響モデル管理サーバが送信した音響モデルにより更新することにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果がある。

【0171】この発明によれば、音響モデル管理サーバが、音声認識装置が音声認識の際に参照する音響モデルを特定するIDを取得し、取得したIDで指示される特定条件に対応して、更新された音響データを読み出し、特定条件に依存した音響モデルを構築して音声認識装置に送信することにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができると共に、特定向けにカスタマイズされた音響モデルを、ネットワークを介して送信することにより、ユーザが複数の異なる音声認識装置を利用する場合でも適切な音響モデルを利用し、高い認識精度を得ることができるという効果がある。

【0172】この発明によれば、言語モデル管理サーバが、更新された言語データを取得して構築した言語モデルを、ネットワークを介して音声認識装置に送信し、音声認識装置が、音声認識の際に参照する言語モデルを、言語モデル管理サーバが送信した言語モデルにより更新することにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果がある。

【0173】この発明によれば、言語モデル管理サーバが、音声認識装置が音声認識の際に参照する言語モデルを特定するIDを取得し、取得したIDで指示される特定条件に対応して、更新された言語データを読み出し、特定条件に依存した言語モデルを構築して音声認識装置

に送信することにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができると共に、特定向けにカスタマイズされた言語モデルを、ネットワークを介して更新することにより、ユーザが複数の異なる音声認識装置を利用する場合でも適切な言語モデルを利用し、高い認識精度を得ることができるという効果がある。

【0174】この発明によれば、音声認識装置が音声認識の際に単語を登録したユーザ辞書を参照し、言語モデル管理サーバが、ネットワークを介してユーザ辞書を読み出し、更新された言語データと、読み出したユーザ辞書とを参照し、ユーザ辞書に依存した言語モデルを構築して音声認識装置に送信することにより、ユーザ辞書が大きくなった場合でも、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果がある。

【0175】この発明によれば、言語モデル管理サーバが、音声認識装置のユーザが利用したテキストを取得し、更新された言語データと、取得したテキストとを参照し、テキストに依存した言語モデルを構築して上記音声認識装置に送信することにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果がある。

【0176】この発明によれば、音声認識装置が、音響モデルを特定するIDと、入力された音声信号から適応化用の音声データとを取得し、取得したID及び適応化用の音声データを、ネットワークを介して音響モデル管理サーバに送信し、音響モデル管理サーバが、送信された適応化用の音声データを用いて、初期音響モデルを適応化し、適応化済み音響モデルを、送信されたIDに対応付けて格納すると共に、外部からの音響モデル更新指令を受けて、ネットワークを介して音声認識装置からIDを受信し、受信したIDに対応する適応化済み音響モデルを、格納している適応化済み音響モデルの中から選択して読み出し、ネットワークを介して音声認識装置に送信し、音声認識装置が、音声認識の際に参照する音響モデルを、音響モデル管理サーバが送信した適応化済み音響モデルにより更新することにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができると共に、ユーザに適応化した適応化済み音響モデルにより、ネットワークを介して音響モデルを更新することで、ユーザが複数の異なる音声認識装置を利用する場合でも、適切な音響モデルを利用し、高い認識精度を得ることができるという効果がある。

【0177】この発明によれば、音声認識装置がネットワークを介して接続された音響モデル管理サーバから、更新された音響データにより構築された音響モデルを受信し、照合手段が音声認識の際に参照する音響モデルを、受信した音響モデルにより更新する音響モデル更新手段とを備えたことにより、ユーザに大きな負担をかけ

ることなく、音声認識の認識精度を向上させることができるという効果がある。

【0178】この発明によれば、音響モデル更新手段が、ネットワークを介して接続された音響モデル管理サーバから、更新された音響データにより構築された、照合手段が音声認識の際に参照する音響モデルの特定条件に依存した音響モデルを受信し、照合手段が音声認識の際に参照する音響モデルを、受信した音響モデルにより更新することにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができると共に、特定向けにカスタマイズされた音響モデルを、ネットワークを介して受信することにより、ユーザが複数の異なる照合手段を利用する場合でも適切な音響モデルを利用し、高い認識精度を得ることができるという効果がある。

【0179】この発明によれば、音声認識装置がネットワークを介して接続された言語モデル管理サーバから、更新された言語データにより構築された言語モデルを受信し、照合手段が音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する言語モデル更新手段とを備えたことにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果がある。

【0180】この発明によれば、言語モデル更新手段が、ネットワークを介して接続された言語モデル管理サーバから、更新された言語データにより構築された、照合手段が音声認識の際に参照する言語モデルの特定条件に依存した言語モデルを受信し、照合手段が音声認識の際に参照する言語モデルを、受信した言語モデルにより更新することにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができると共に、特定向けにカスタマイズされた言語モデルを、ネットワークを介して受信することにより、ユーザが複数の異なる照合手段を利用する場合でも適切な言語モデルを利用し、高い認識精度を得ることができるという効果がある。

【0181】この発明によれば、照合手段が音声認識の際に参照する単語を登録したユーザ辞書を備え、言語モデル更新手段が、ネットワークを介して接続された言語モデル管理サーバから、更新された言語データにより構築された、ユーザ辞書に依存した言語モデルを受信し、照合手段が音声認識の際に参照する言語モデルを、受信した言語モデルにより更新することにより、ユーザ辞書が大きくなった場合でも、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果がある。

【0182】この発明によれば、言語モデル更新手段が、ネットワークを介して接続された言語モデル管理サーバから、更新された言語データにより構築された、音声認識を行うユーザが利用したテキストに依存した言語

モデルを受信し、照合手段が音声認識の際に参照する言語モデルを、受信した言語モデルにより更新することにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果がある。

【0183】この発明によれば、音声の音響的な観測値系列の確率を求める音響モデルと、音声信号を入力し上記音響モデルを参照して音声認識を行い、認識結果を出力する照合手段と、音響モデルを特定するIDを取得する音響モデルID取得手段と、取得したIDを読み出し、入力された音声信号から適応化用の音声データを取得し、読み出したID及び取得した適応化用の音声データを、ネットワークを介して接続された音響モデル管理サーバに送信する適応化用音声取得手段と、音響モデル管理サーバから、IDに対応する適応化用の音声データにより適応化された適応化済み音響モデルを受信し、照合手段が音声認識の際に参照する音響モデルを、受信した適応化済み音響モデルにより更新する音響モデル更新手段とを備えたことにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができると共に、ユーザに適応化した適応化済み音響モデルにより、ネットワークを介して音響モデルを更新することで、ユーザが複数の異なる照合手段を利用する場合でも、適切な音響モデルを利用し、高い認識精度を得ることができるという効果がある。

【0184】この発明によれば、更新された音響データを取得する音響データ取得手段と、外部からの音響モデル更新指令を受けて更新された音響データを読み出し、音声の音響的な観測値系列の確率を求める音響モデルを構築する音響モデル構築手段と、音響モデル構築手段により構築された音響モデルを、ネットワークを介して音声認識を行う音声認識装置に送信する音響モデル送信手段とを備えたことにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果がある。

【0185】この発明によれば、更新された音響データを取得する音響データ取得手段と、外部からの音響モデル更新指令を受けて、音声認識の際に参照する音響モデルを特定するIDを取得する更新音響モデルID取得手段と、取得したIDで指示される特定条件に対応して、更新された音響データを読み出す特定向け音響データ読み出し手段と、読み出した更新された音響データを参照し、特定条件に依存した音響モデルを構築する特定向け音響モデル構築手段と、構築した音響モデルを、ネットワークを介して音声認識装置に送信する音響モデル送信手段とを備えたことにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができると共に、特定向けにカスタマイズされた音響モデルを、ネットワークを介して送信することにより、ユーザが複数の異なる音声認識装置を利用する場合でも適切な

音響モデルを利用し、高い認識精度を得ることができるという効果がある。

【0186】この発明によれば、音声の音響的な観測値系列の確率を求める、適応化前の初期音響モデルと、ネットワークを介して接続された音声認識装置から送信された、適応化用の音声データと、音声認識装置が音声認識の際に参照する音響モデルを特定するIDを受信し、適応化用の音声データを用いて初期音響モデルを適応化し、適応化済み音響モデルを、受信したIDに対応付けて適応化済み音響モデル格納手段に格納する音響モデル適応化手段と、外部からの音響モデル更新指令を受けて、ネットワークを介して音声認識装置からIDを受信し、受信したIDに対応する適応化済み音響モデルを、適応化済み音響モデル格納手段から選択して読み出す適応化済み音響モデル選択手段と、読み出した適応化済み音響モデルを、ネットワークを介して音声認識装置に送信する音響モデル送信手段とを備えたことにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができると共に、ユーザに適応化した適応化済み音響モデルにより、ネットワークを介して音響モデルを更新することで、ユーザが複数の異なる音声認識装置を利用する場合でも、適切な音響モデルを利用し、高い認識精度を得ることができるという効果がある。

【0187】この発明によれば、更新された言語データを取得する言語データ取得手段と、外部からの言語モデル更新指令を受けて更新された言語データを読み出し、単語列の出現確率を求める言語モデルを構築する言語モデル構築手段と、言語モデル構築手段が構築した言語モデルを、ネットワークを介して音声認識を行う音声認識装置に送信する言語モデル送信手段とを備えたことにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果がある。

【0188】この発明によれば、更新された言語データを取得する言語データ取得手段と、外部からの言語モデル更新指令を受けて、音声認識の際に参照する言語モデルを特定するIDを取得する更新言語モデルID取得手段と、IDで指示される特定条件に対応して、更新された言語データを読み出す特定向け言語データ読み出し手段と、読み出した更新された言語データを参照し、特定条件に依存した言語モデルを構築する特定向け言語モデル構築手段と、構築した言語モデルを、ネットワークを介して音声認識装置に送信する言語モデル送信手段とを備えたことにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができると共に、特定向けにカスタマイズされた言語モデルを、ネットワークを介して送信することにより、ユーザが複数の異なる音声認識装置を利用する場合でも適切な言語モデルを利用し、高い認識精度を得ることができるという効果がある。

【0189】この発明によれば、更新された言語データを取得する言語データ取得手段と、外部からの言語モデル更新指令を受けて、ネットワークを介して接続された音声認識装置が音声認識の際に参照するユーザ辞書を読み出すユーザ辞書読み出し手段と、更新された言語データを読み出し、ユーザ辞書に依存した言語モデルを構築するユーザ辞書依存言語モデル構築手段と、構築した言語モデルを、ネットワークを介して音声認識装置に送信する言語モデル送信手段とを備えたことにより、ユーザ辞書が大きくなった場合でも、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果がある。

【0190】この発明によれば、更新された言語データを取得する言語データ取得手段と、外部からの言語モデル更新指令を受けて、ユーザが利用したテキストを取得するユーザ利用テキスト取得手段と、更新された言語データを読み出し、取得したテキストに依存した言語モデルを構築するユーザ利用テキスト依存言語モデル構築手段と、構築した言語モデルを、ネットワークを介して音声認識装置に送信する言語モデル送信手段とを備えたことにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果がある。

【0191】この発明によれば、更新された音響データを取得する第1のステップと、音響モデル更新指令を受けて更新された音響データを読み出し、音響モデルを構築する第2のステップと、構築した音響モデルを、ネットワークを介して送信する第3のステップと、音響モデルを受信し、音声認識の際に参照する音響モデルを、受信した音響モデルにより更新する第4のステップとを備えたことにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果がある。

【0192】この発明によれば、更新された言語データを取得する第1のステップと、言語モデル更新指令を受けて更新された言語データを読み出し、言語モデルを構築する第2のステップと、構築した言語モデルを、ネットワークを介して送信する第3のステップと、言語モデルを受信し、音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する第4のステップとを備えたことにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果がある。

【0193】この発明によれば、更新された音響データを取得する第1のステップと、音響モデル更新指令を受けて、音声認識の際に参照する音響モデルを特定するIDを取得する第2のステップと、IDで指示される特定条件に対応して、更新された音響データを読み出す第3のステップと、更新された音響データを参照し、特定条件に依存した音響モデルを構築する第4のステップと、

構築した音響モデルを、ネットワークを介して送信する第5のステップと、音響モデルを受信し、音声認識の際に参照する音響モデルを、受信した音響モデルにより更新する第6のステップとを備えたことにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができると共に、特定向けにカスタマイズされた音響モデルを、ネットワークを介して受信することにより、ユーザが複数の異なる音声認識方法を利用する場合でも適切な音響モデルを利用し、高い認識精度を得ることができるという効果がある。

【0194】この発明によれば、更新された言語データを取得する第1のステップと、言語モデル更新指令を受けて、音声認識の際に参照する言語モデルを特定するIDを取得する第2のステップと、取得したIDで指示される特定条件に対応して、更新された言語データを読み出す第3のステップと、更新された言語データを参照し、特定条件に依存した言語モデルを構築する第4のステップと、構築した言語モデルを、ネットワークを介して送信する第5のステップと、言語モデルを受信し、音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する第6のステップとを備えたことにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができると共に、特定向けにカスタマイズされた言語モデルを、ネットワークを介して受信することにより、ユーザが複数の異なる音声認識方法を利用する場合でも適切な言語モデルを利用し、高い認識精度を得ることができるという効果がある。

【0195】この発明によれば、更新された言語データを取得する第1のステップと、言語モデル更新指令を受けて、音声認識の際に参照するユーザ辞書を読み出す第2のステップと、更新された言語データを読み出し、ユーザ辞書に依存した言語モデルを構築する第3のステップと、構築した言語モデルを、ネットワークを介して送信する第4のステップと、言語モデルを受信し、音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する第5のステップとを備えたことにより、ユーザ辞書が大きくなった場合でも、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果がある。

【0196】この発明によれば、更新された言語データを取得する第1のステップと、言語モデル更新指令を受けて、ユーザが利用したテキストを取得する第2のステップと、更新された言語データを読み出し、テキストに依存した言語モデルを構築する第3のステップと、構築した言語モデルを、ネットワークを介して送信する第4のステップと、言語モデルを受信し、音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する第5のステップとを備えたことにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果がある。

【0197】この発明によれば、音響モデルを特定するIDを取得する第1のステップと、取得したIDを読み出し、入力された音声信号から適応化用の音声データを取得し、ネットワークを介して、読み出したID及び取得した適応化用の音声データを送信する第2のステップと、送信した適応化用の音声データを用いて、適応化前の初期音響モデルを適応化し、適応化済み音響モデルを、送信したIDに対応付けて格納する第3のステップと、音響モデル更新指令を受けて、ネットワークを介して第1のステップで取得したIDを受信し、受信したIDに対応する適応化済み音響モデルを、第3のステップで格納している適応化済み音響モデルの中から選択して読み出す第4のステップと、第4のステップで読み出した適応化済み音響モデルを、ネットワークを介して送信する第5のステップと、送信した適応化済み音響モデルを受信し、音声認識の際に参照する音響モデルを、受信した適応化済み音響モデルにより更新する第6のステップとを備えたことにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができると共に、ユーザに適応化した適応化済み音響モデルにより、ネットワークを介して音響モデルを更新することで、ユーザが複数の異なる音声認識方法を利用する場合でも、適切な音響モデルを利用し、高い認識精度を得ることができるという効果がある。

【0198】この発明によれば、音声認識プログラムを記録した記録媒体で、更新された音響データを取得する音響データ取得機能と、音響モデル更新指令を受けて更新された音響データを読み出し、音響モデルを構築する音響モデル構築機能と、構築した音響モデルを、ネットワークを介して送信する音響モデル送信機能と、音響モデルを受信し、照合機能が音声認識の際に参照する音響モデルを、受信した音響モデルにより更新する音響モデル更新機能とを実現させることにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果がある。

【0199】この発明によれば、音声認識プログラムを記録した記録媒体で、更新された言語データを取得する言語データ取得機能と、言語モデル更新指令を受けて、上記言語データ取得機能が取得した更新された言語データを読み出し、言語モデルを構築する言語モデル構築機能と、上記言語モデル構築機能が構築した言語モデルを、ネットワークを介して送信する言語モデル送信機能と、上記言語モデル送信機能が送信した言語モデルを受信し、上記照合機能が音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する言語モデル更新機能とを実現させることにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果がある。

【0200】この発明によれば、音声認識プログラムを記録した記録媒体で、更新された音響データを取得する



音響データ取得機能と、音響モデル更新指令を受けて、音響モデルを特定するIDを取得する更新音響モデルID取得機能と、取得したIDで指示される特定条件に対応して、更新された音響データを読み出す特定向け音響データ読み出し機能と、更新された音響データを参照し、特定条件に依存した音響モデルを構築する特定向け音響モデル構築機能と、音響モデルを、ネットワークを介して送信する音響モデル送信機能と、音響モデルを受信し、照合機能が音声認識の際に参照する音響モデルを、受信した音響モデルにより更新する音響モデル更新機能とを実現させることにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができると共に、特定向けにカスタマイズされた音響モデルを、ネットワークを介して受信することにより、ユーザが複数の照合機能を利用する場合でも適切な音響モデルを利用し、高い認識精度を得ることができるという効果がある。

【0201】この発明によれば、音声認識プログラムを記録した記録媒体で、更新された言語データを取得する言語データ取得機能と、言語モデル更新指令を受けて、言語モデルを特定するIDを取得する更新言語モデルID取得機能と、取得したIDで指示される特定条件に対応して、更新された言語データを読み出す特定向け言語データ読み出し機能と、更新された言語データを参照し、特定条件に依存した言語モデルを構築する特定向け言語モデル構築機能と、構築した言語モデルを、ネットワークを介して送信する言語モデル送信機能と、言語モデルを受信し、照合機能が音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する言語モデル更新機能とを実現させることにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができると共に、特定向けにカスタマイズされた言語モデルを、ネットワークを介して受信することにより、ユーザが複数の照合機能を利用する場合でも適切な言語モデルを利用し、高い認識精度を得ることができるという効果がある。

【0202】この発明によれば、音声認識プログラムを記録した記録媒体で、更新された言語データを取得する言語データ取得機能と、言語モデル更新指令を受けて、ユーザ辞書を読み出すユーザ辞書読み出し機能と、更新された言語データを読み出し、ユーザ辞書に依存した言語モデルを構築するユーザ辞書依存言語モデル構築機能と、構築した言語モデルを、ネットワークを介して送信する言語モデル送信機能と、言語モデルを受信し、照合機能が音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する言語モデル更新機能とを実現させることにより、ユーザ辞書が大きくなった場合でも、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果がある。

【0203】この発明によれば、音声認識プログラムを

記録した記録媒体で、更新された言語データを取得する言語データ取得機能と、言語モデル更新指令を受けて、ユーザが利用したテキストを取得するユーザ利用テキスト取得機能と、更新された言語データを読み出し、テキストに依存した言語モデルを構築するユーザ利用テキスト依存言語モデル構築機能と、構築した言語モデルを、ネットワークを介して送信する言語モデル送信機能と、言語モデルを受信し、照合機能が音声認識の際に参照する言語モデルを、受信した言語モデルにより更新する言語モデル更新機能とを実現させることにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができるという効果がある。

【0204】この発明によれば、音声認識プログラムを記録した記録媒体で、音響モデルを特定するIDを取得する音響モデルID取得機能と、取得したIDを読み出し、入力された音声信号から適応化用の音声データを取得し、ネットワークを介して、読み出したID及び取得した適応化用の音声データを送信する適応化用音声取得機能と、送信した適応化用の音声データを用いて、適応化前の初期音響モデルを適応化し、適応化済み音響モデルを、送信したIDに対応付けて格納する音響モデル適応化機能と、音響モデル更新指令を受けて、ネットワークを介して音響モデルID取得機能が取得したIDを受信し、受信したIDに対応する適応化済み音響モデルを、音響モデル適応化機能が格納した適応化済み音響モデルの中から選択して読み出す適応化済み音響モデル選択機能と、読み出した適応化済み音響モデルを、ネットワークを介して送信する音響モデル送信機能と、適応化済み音響モデルを受信し、照合機能が音声認識の際に参照する音響モデルを、受信した適応化済み音響モデルにより更新する音響モデル更新機能とを実現させることにより、ユーザに大きな負担をかけることなく、音声認識の認識精度を向上させることができると共に、ユーザに適応化した適応化済み音響モデルにより、ネットワークを介して音響モデルを更新することで、ユーザが複数の異なる照合機能を利用する場合でも、適切な音響モデルを利用し、高い認識精度を得ることができるという効果がある。

【図面の簡単な説明】

【図1】 この発明の実施の形態1による音声認識装置の構成を示すブロック図である。

【図2】 この発明の実施の形態1による音響モデルの更新処理の手順を示すフローチャートである。

【図3】 この発明の実施の形態1による言語モデルの更新処理の手順を示すフローチャートである。

【図4】 この発明の実施の形態2による音声認識装置の構成を示すブロック図である。

【図5】 この発明の実施の形態2による言語モデルの更新処理の手順を示すフローチャートである。

【図6】 この発明の実施の形態3による音声認識装置

の構成を示すブロック図である。

【図7】 この発明の実施の形態3による言語モデルの更新処理の手順を示すフローチャートである。

【図8】 この発明の実施の形態4による音声認識装置の構成を示すブロック図である。

【図9】 この発明の実施の形態4による言語モデルの更新処理の手順を示すフローチャートである。

【図10】 この発明の実施の形態5による音声認識装置の構成を示すブロック図である。

【図11】 この発明の実施の形態5による音響モデルの更新処理の手順を示すフローチャートである。

【図12】 従来の音声認識装置の構成を示すブロック図である。

【図13】 従来の音声認識処理の手順を示すフローチャートである。

【図14】 従来の音声認識装置の構成を示すブロック図である。

【図15】 従来の音声認識装置の構成を示すブロック図である。

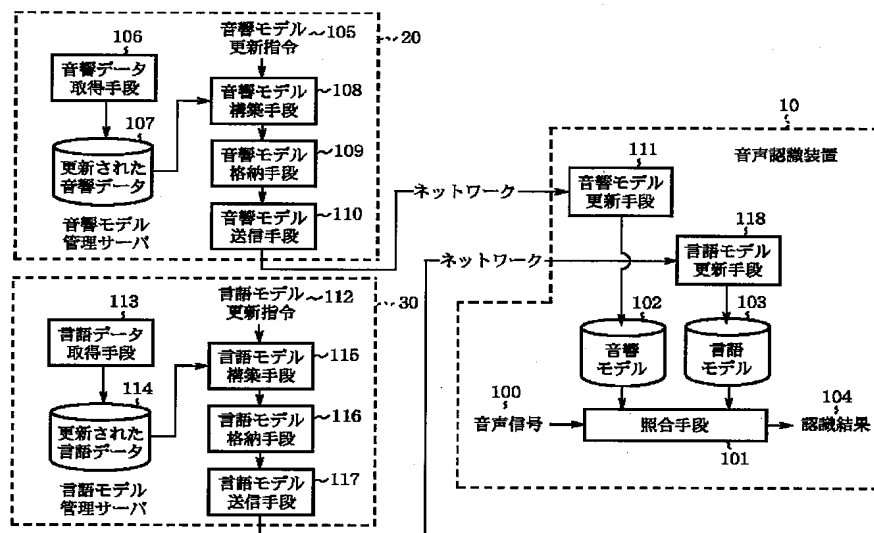
【図16】 ユーザ辞書601の構成例を示す図である。

【図17】 従来の音声認識処理の手順を示すフローチャートである。

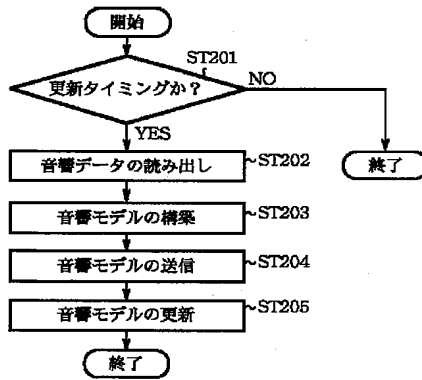
＊【符号の説明】

10 音声認識装置、20 音響モデル管理サーバ、30 言語モデル管理サーバ、100 音声信号、101 照合手段、102 音響モデル、103 言語モデル、104 認識結果、105 音響モデル更新指令、106 音響データ取得手段、107 更新された音響データ、108 音響モデル構築手段、109 音響モデル格納手段、110 音響モデル送信手段、111 音響モデル更新手段、112 言語モデル更新指令、113 言語データ取得手段、114 更新された言語データ、115 言語モデル構築手段、116 言語モデル格納手段、117 言語モデル送信手段、401 更新言語モデルID取得手段、402 特定向け言語データ読み出し手段、403 特定向け言語モデル構築手段、601 ユーザ辞書、602 ユーザ辞書読み出し手段、603 ユーザ辞書依存言語モデル構築手段、801 ユーザ利用テキスト取得手段、802 ユーザ利用テキスト格納手段、803 ユーザ利用テキスト依存言語モデル構築手段、1001 音響モデルID取得手段、1002 適応化用音声取得手段、1003 初期音響モデル、1004 音響モデル適応化手段、1005 適応化済み音響モデル格納手段、1006 適応化済み音響モデル選択手段。

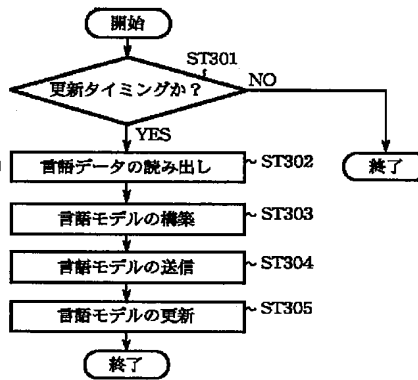
【図1】



【図2】



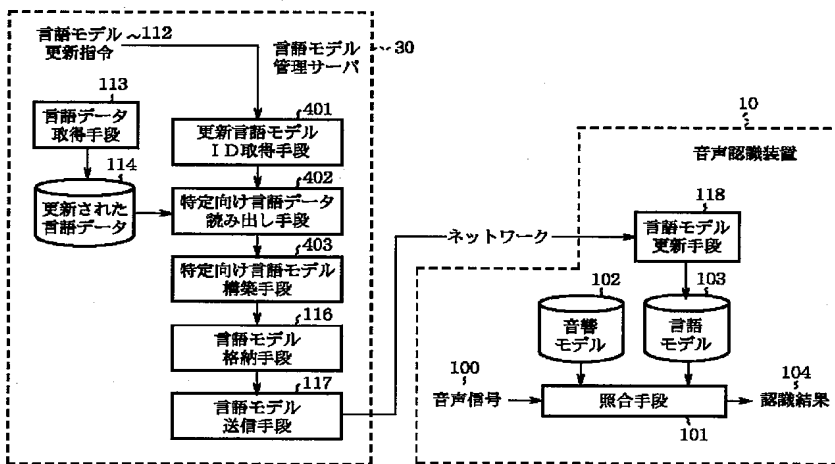
【図3】



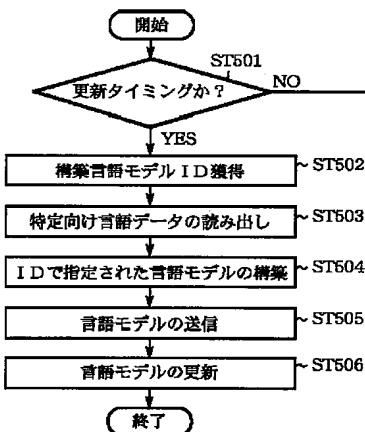
【図16】

漢字表記	よみ
音素	オンソ
音声認識	オンセイニンシキ
多変量解析	タヘンリョウカイセキ
チューリング	チューリング
...	...

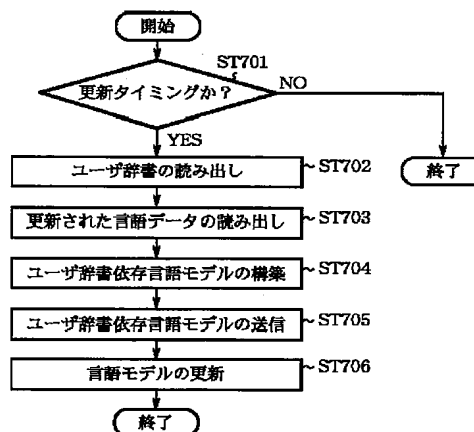
【図4】



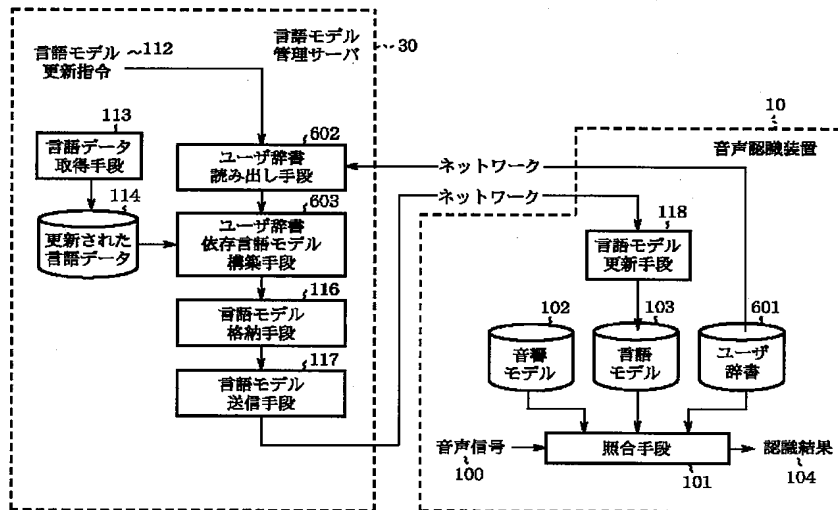
【図5】



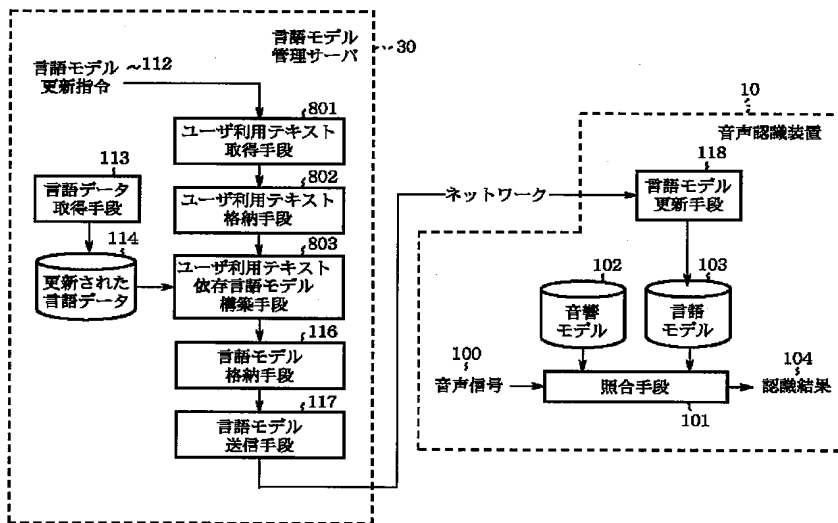
【図7】



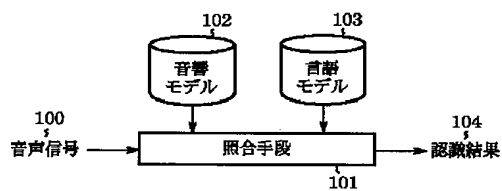
【図6】



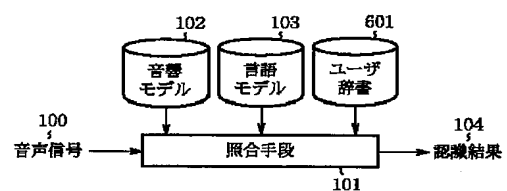
【図8】



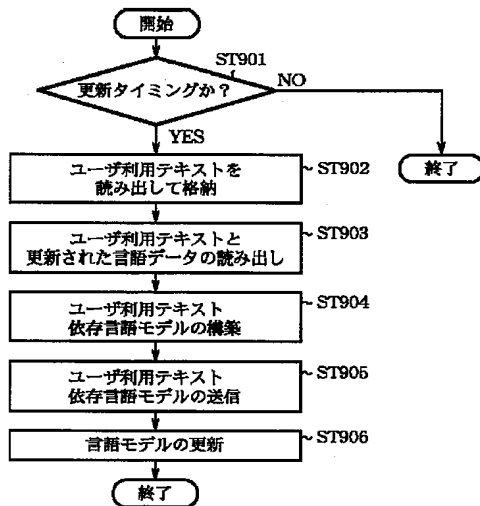
【図12】



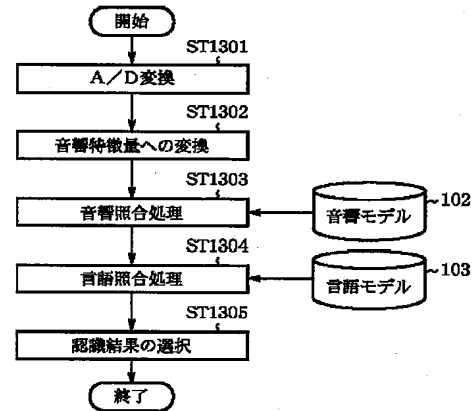
【図15】



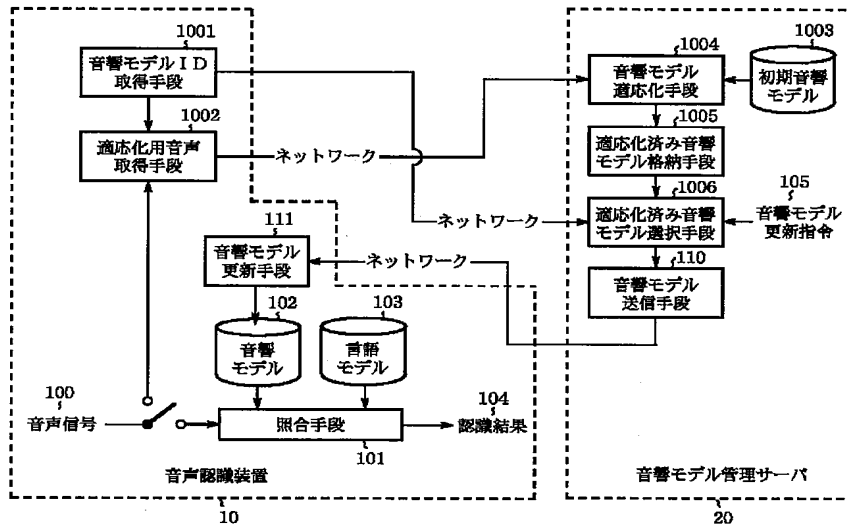
【図9】



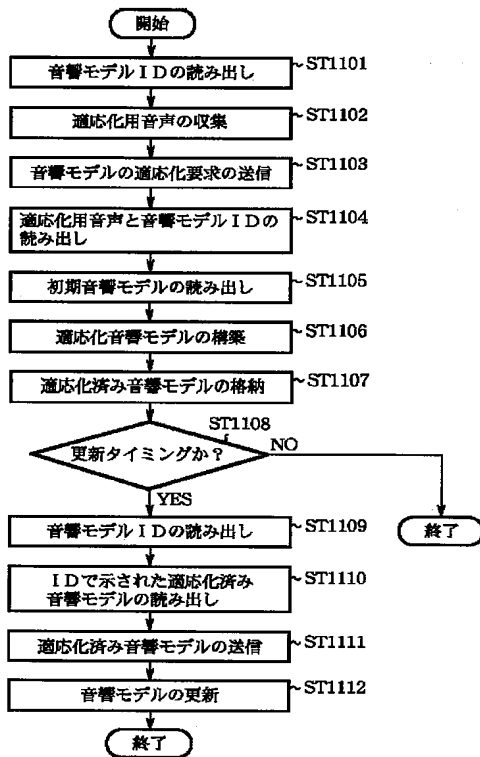
【図13】



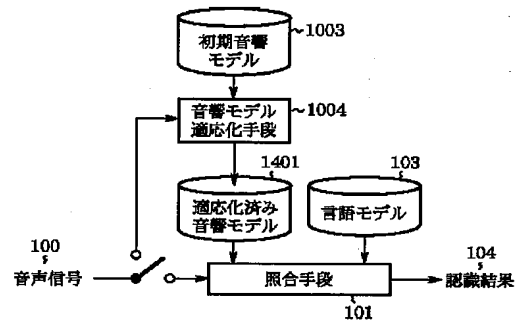
【図10】



【図11】



【図14】



【図17】

